# Provenance-Enabled Data Exploration and Visualization Tutorial

Erik Anderson
Juliana Freire
David Koop
Emanuele Santos
Claudio Silva
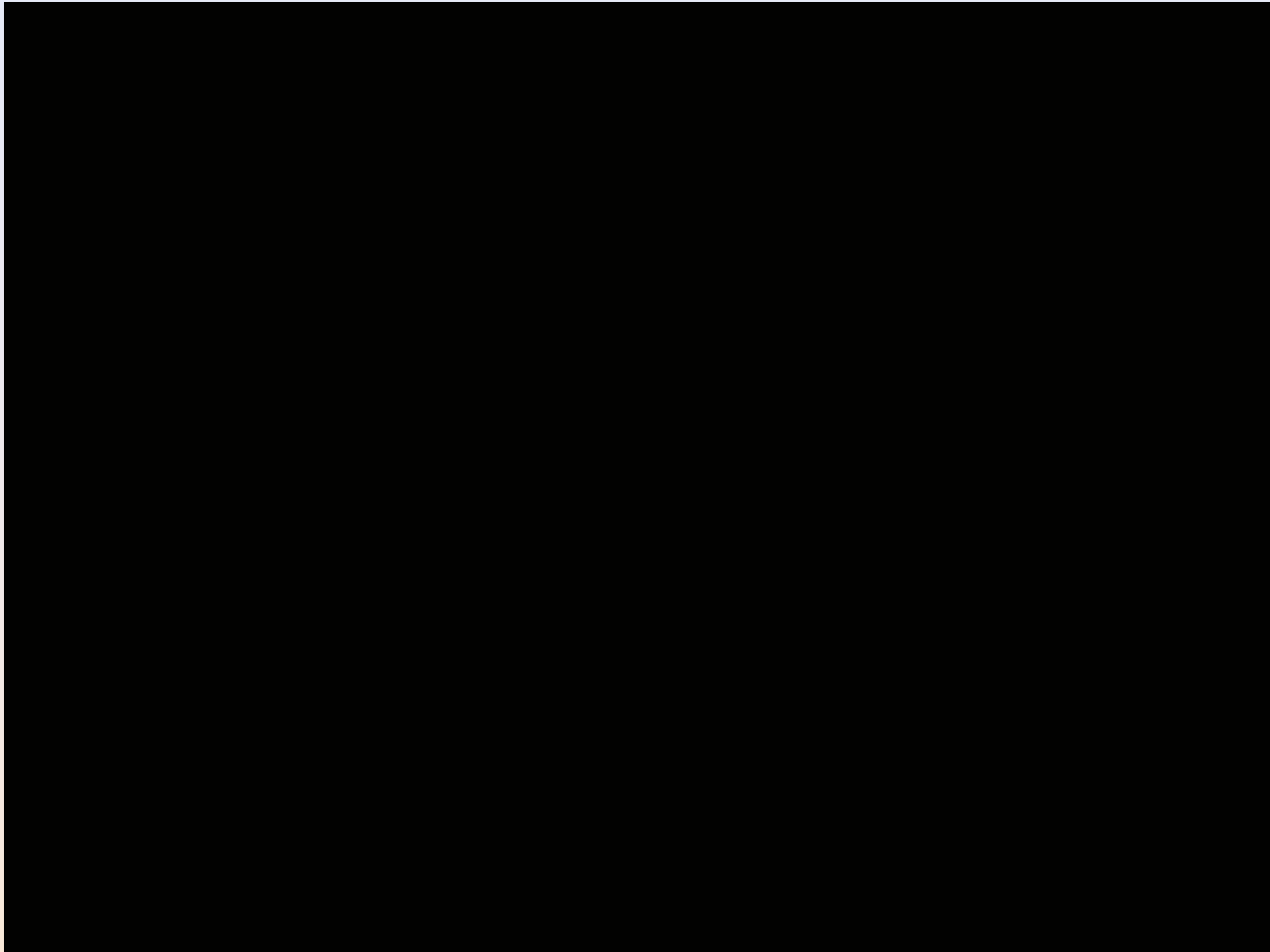
# Presentation VIS

"Study of a Numerically Modeled Severe Storm", NCSA, UIUC

# Presentation VIS



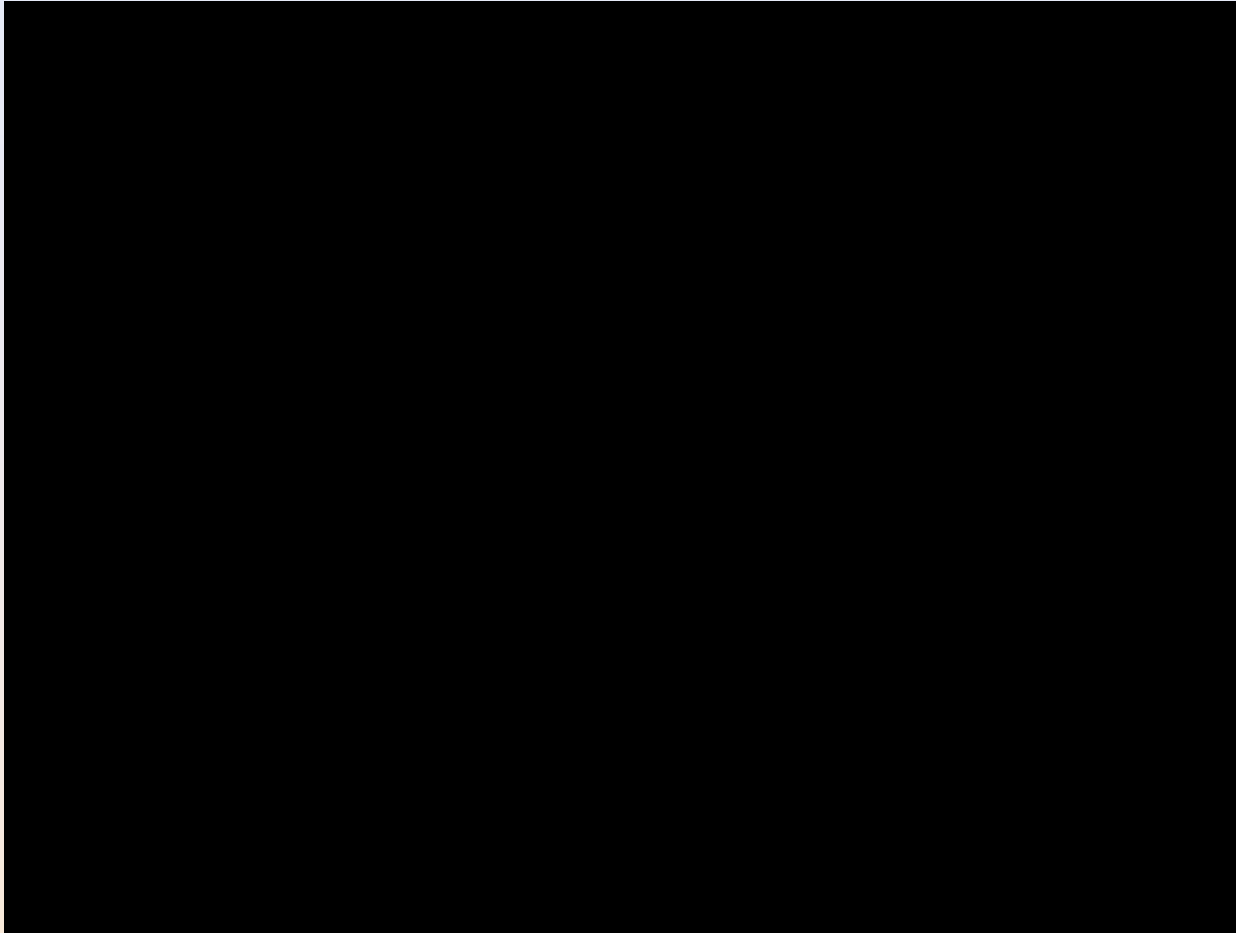"Study of a Numerically Modeled Severe Storm", NCSA, UIUC

# Presentation VIS

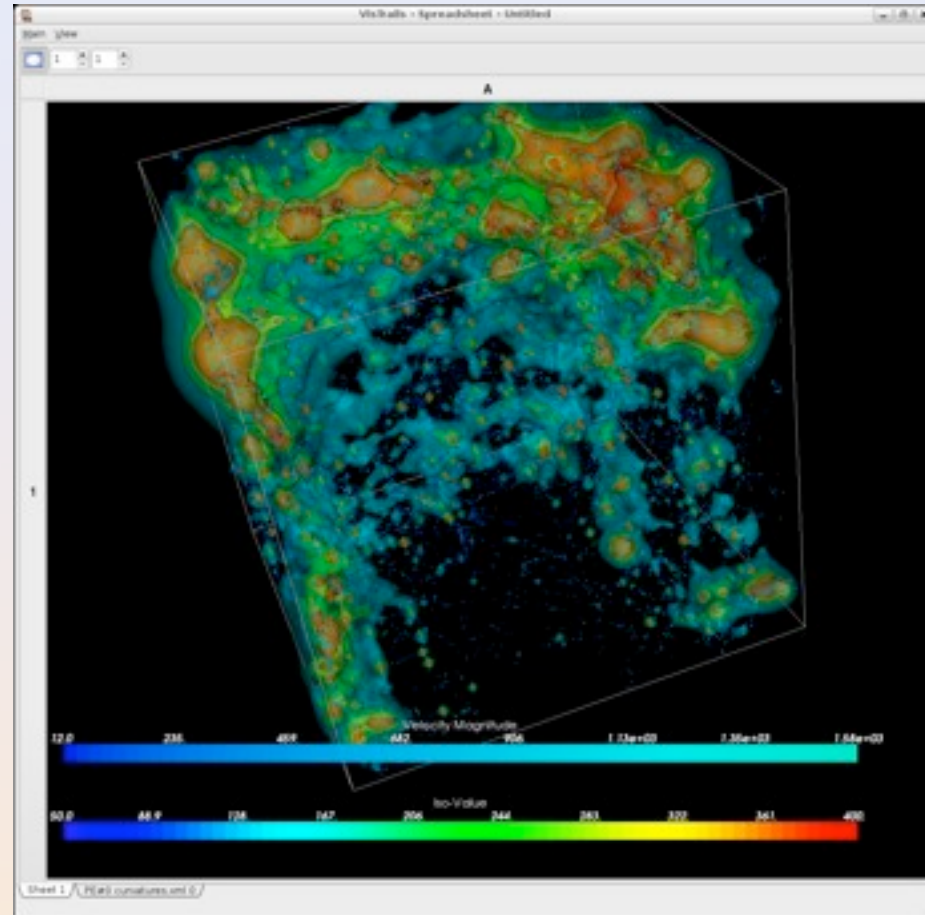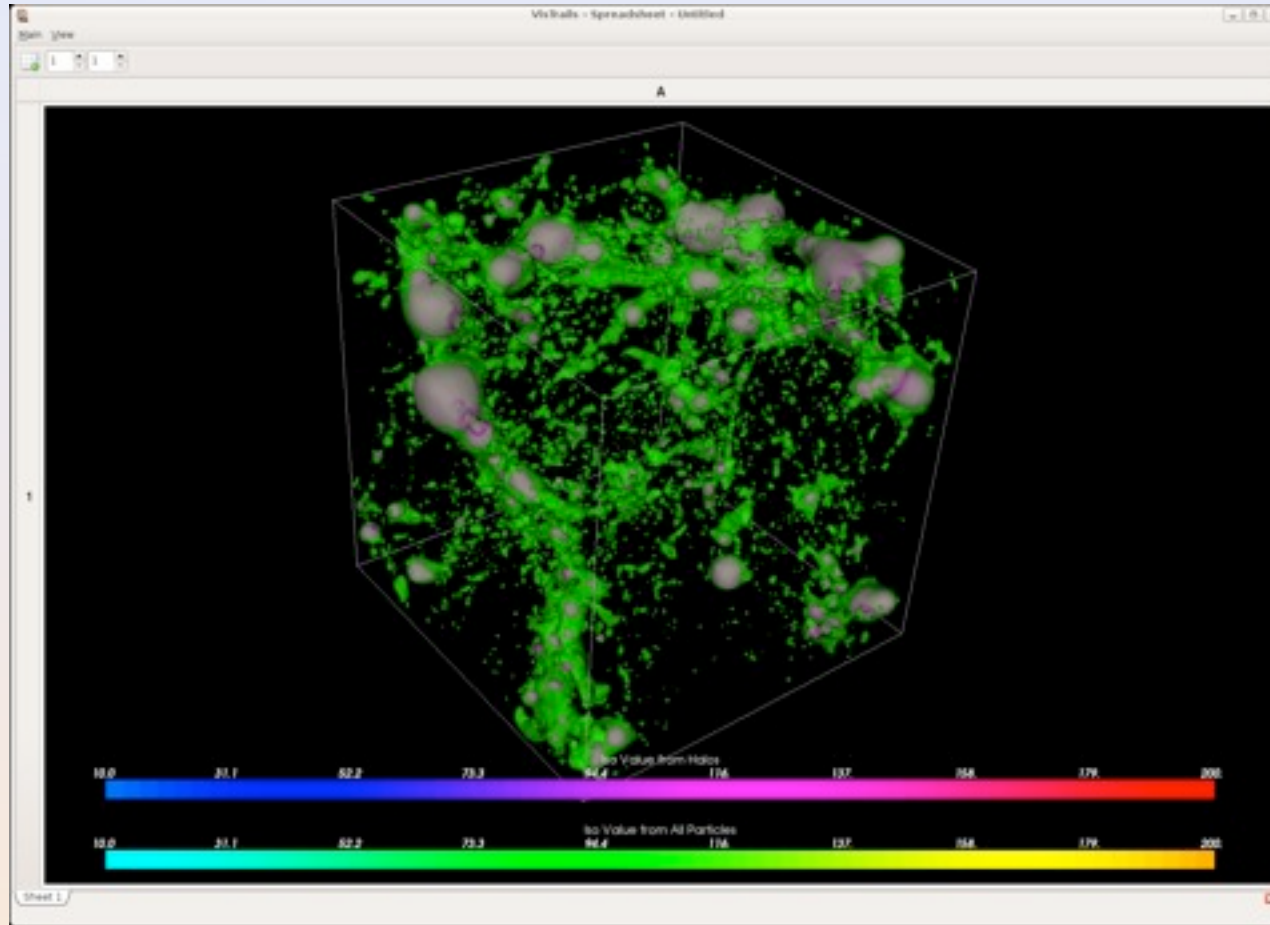"Fusion Simulation Visualization" Kruger, Sanderson, et al

# Presentation VIS



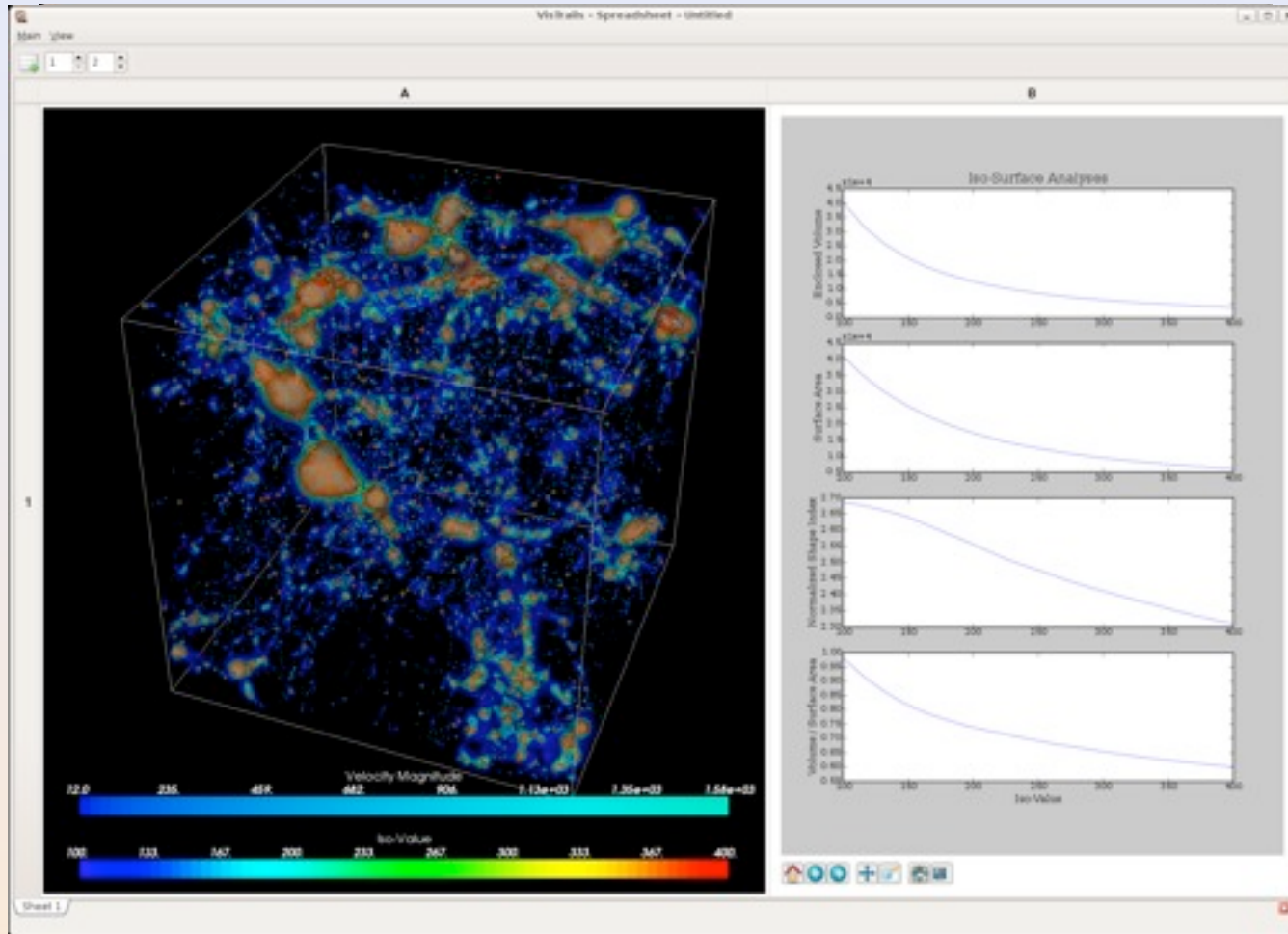"Fusion Simulation Visualization" Kruger, Sanderson, et al

# Exploratory VIS



"The Cosmic Code Comparison Project," Ahrens, Anderson, Heitmann, Habib, et al
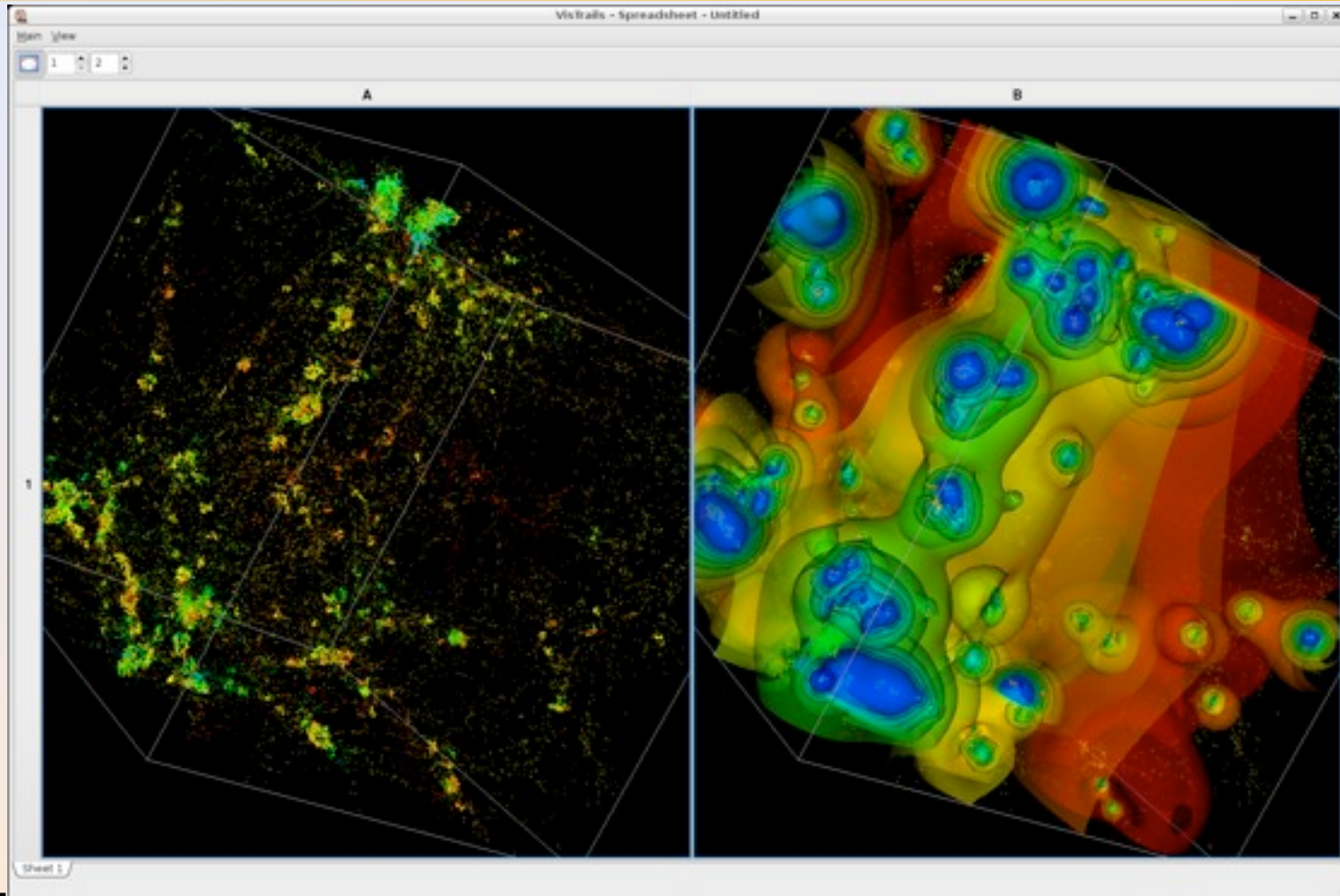
# Exploratory VIS



"The Cosmic Code Comparison Project," Ahrens, Anderson, Heitmann, Habib, et al

# Exploratory VIS



"The Cosmic Code Comparison Project," Ahrens, Anderson, Heitmann, Habib, et al

# Exploratory VIS



"The Cosmic Code Comparison Project," Ahrens, Anderson, Heitmann, Habib, et al
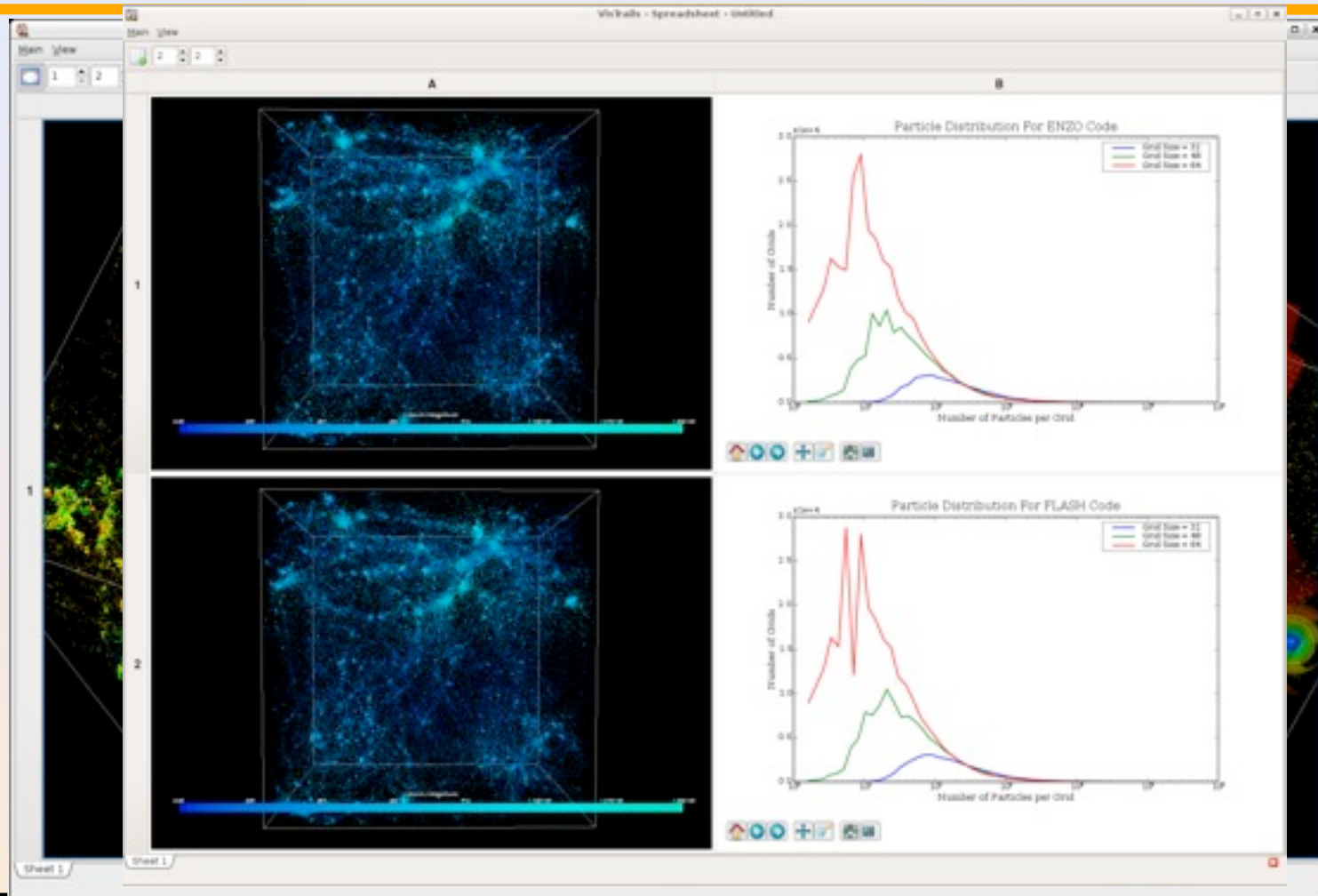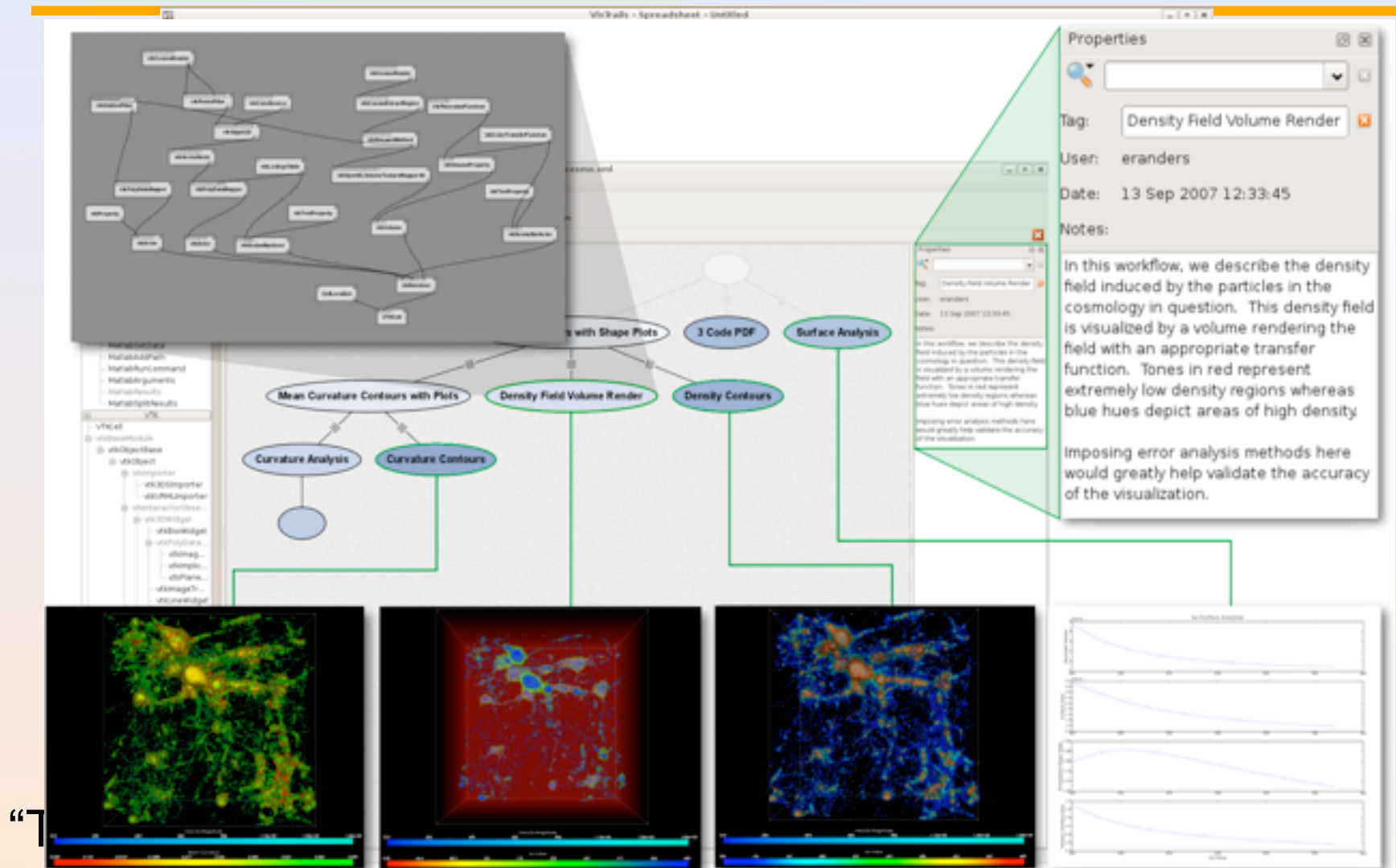
# Exploratory VIS



"The Cosmic Code Comparison Project," Ahrens, Anderson, Heitmann, Habib, et al
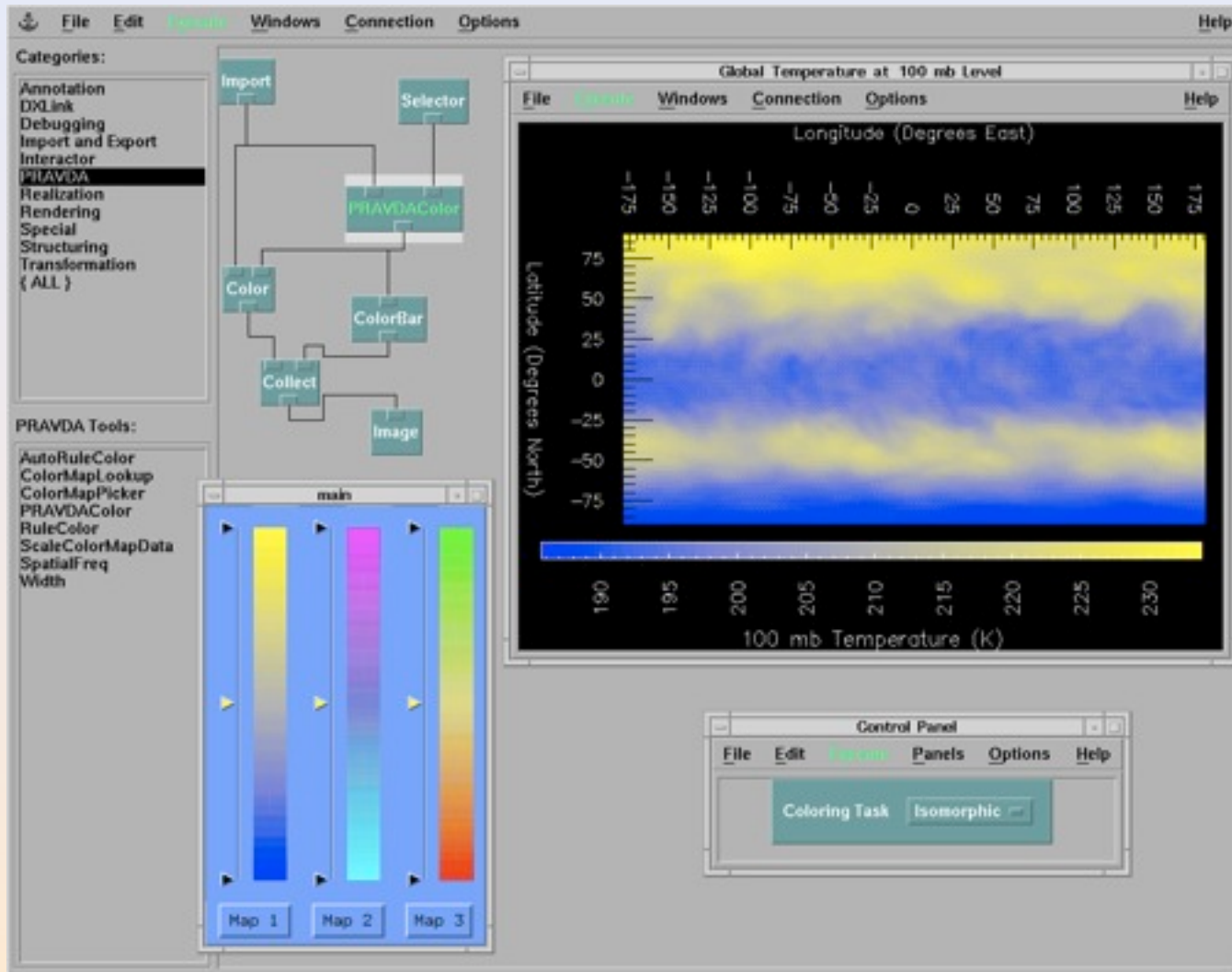
# Exploratory VIS



"T
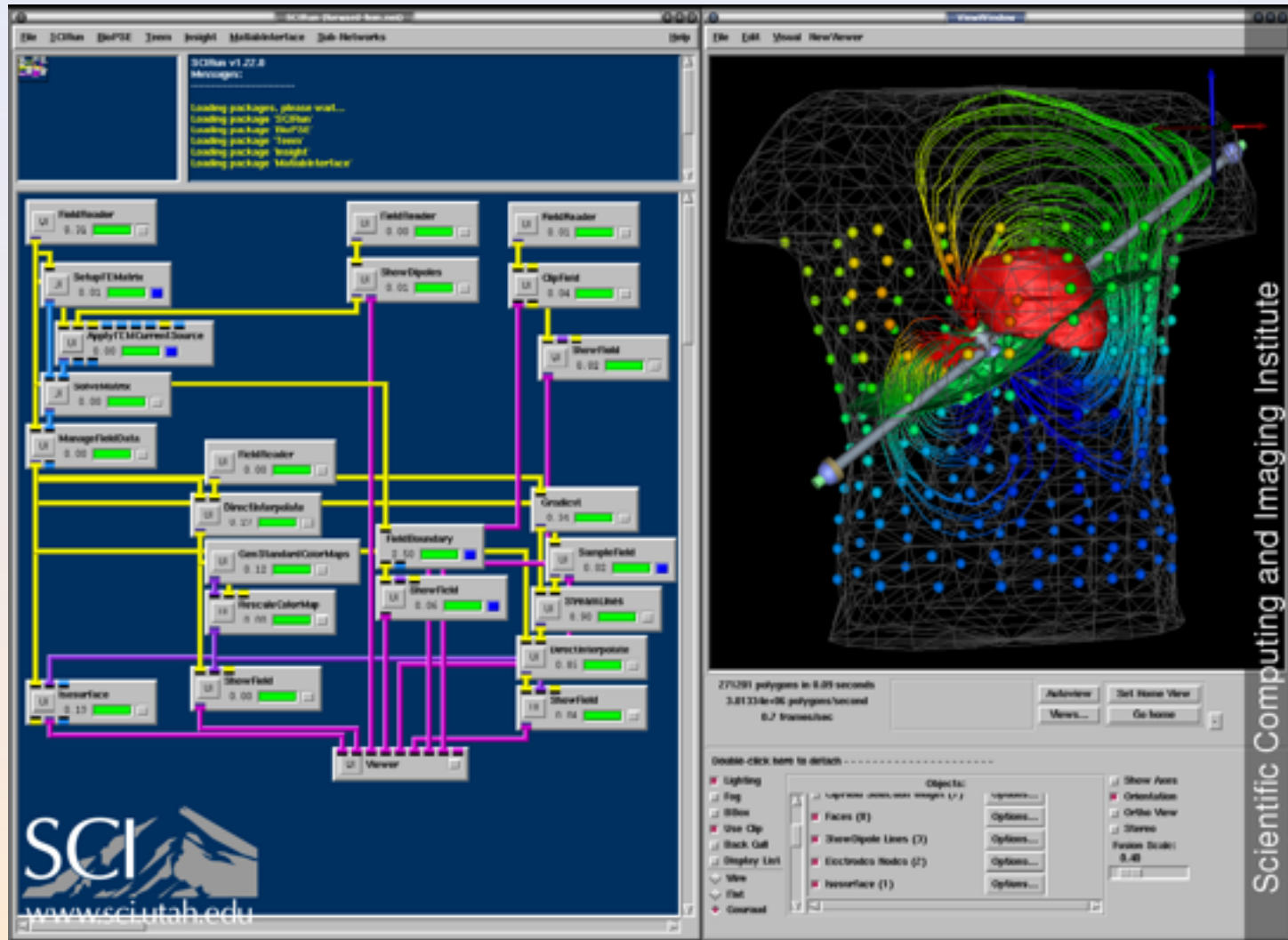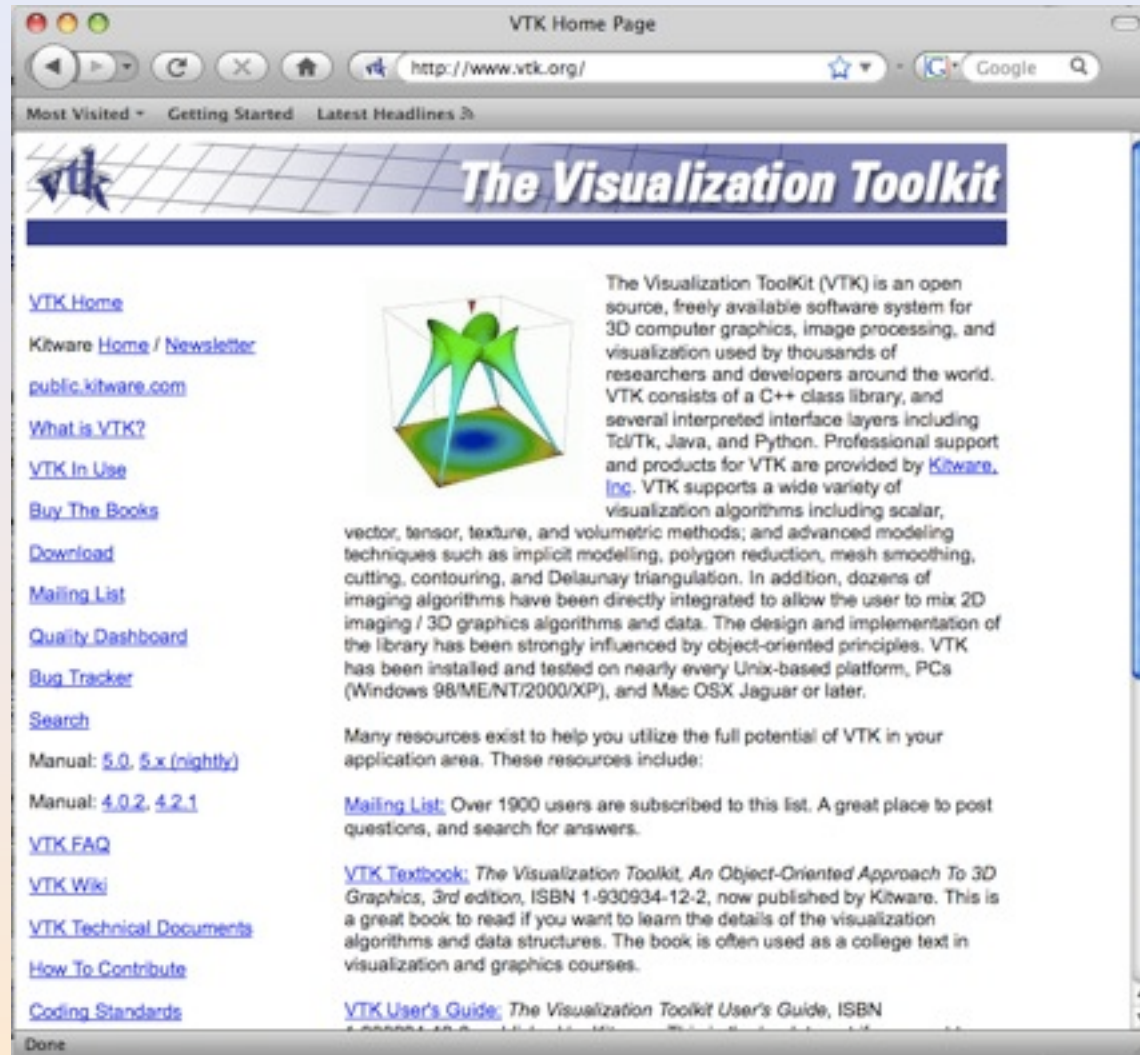Anderson, Heitmann, Habib, et al

# Influential VIS tools

# IBM OpenDX

# SCIRun

# VTK

# Digression: Workflows
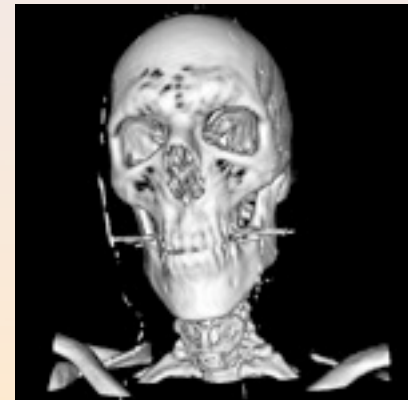
# Scientific Workflows and Dataflows

Dataflows are directed graphs describing a computational task

- Vertices = modules = processing steps + **parameters**
- Edges = connections between output and input ports
- Execution order determined by flow of data from output to input ports
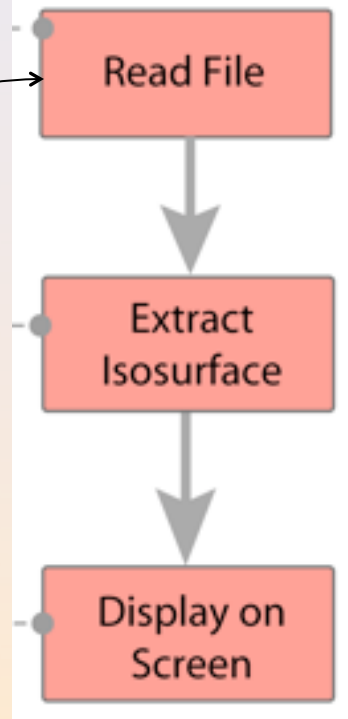
Input:
Head.120.iso

Output:

# Scientific Workflows and Dataflows

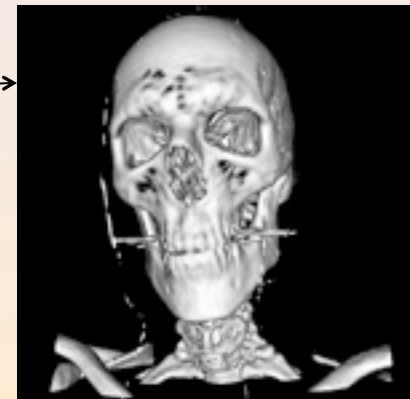Dataflows are directed graphs describing a computational task
- Vertices = modules = processing steps + **parameters**
- Edges = connections between output and input ports
- Execution order determined by flow of data from output to input ports

Input:
Head.120.iso

Read File

Extract Isosurface

Display on Screen

**Isosurface**
*value=57*

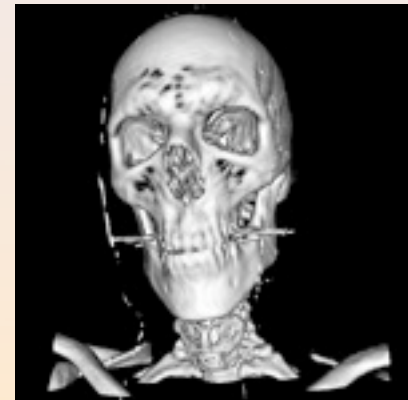Output:

# Scientific Workflows and Dataflows

- Dataflows are directed graphs describing a computational task
  - Vertices = modules = processing steps + **parameters**
  - Edges = connections between output and input ports
  - Execution order determined by flow of data from output to input ports

```
Input:
Head.120.iso
```

Output:

# Scientific Workflows and Dataflows

- Dataflows are directed graphs describing a computational task
  - Vertices = modules = processing steps + **parameters**
  - Edges = connections between output and input ports
  - Execution order determined by flow of data from output to input ports
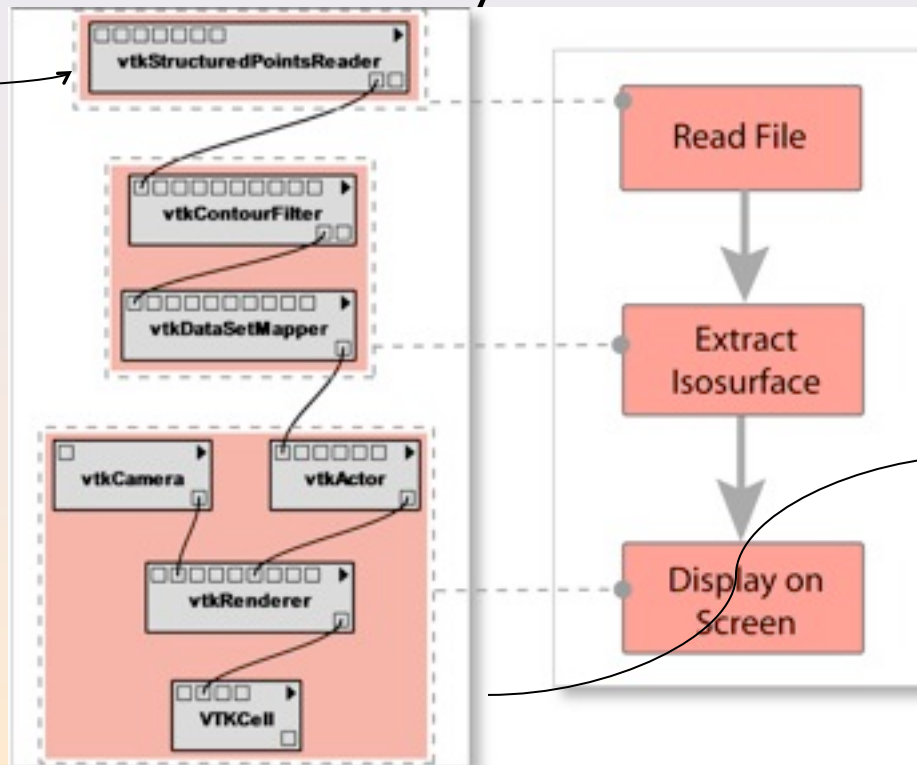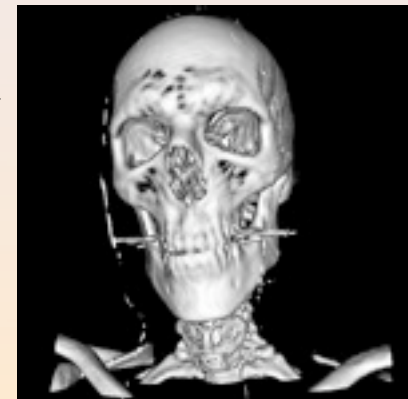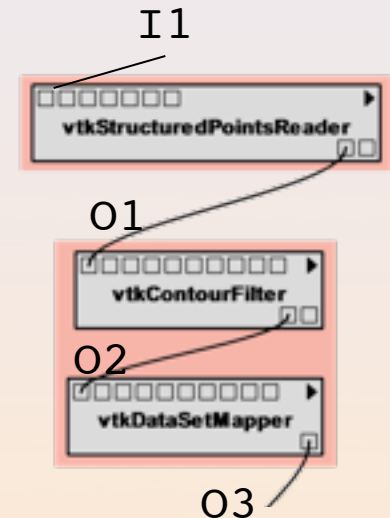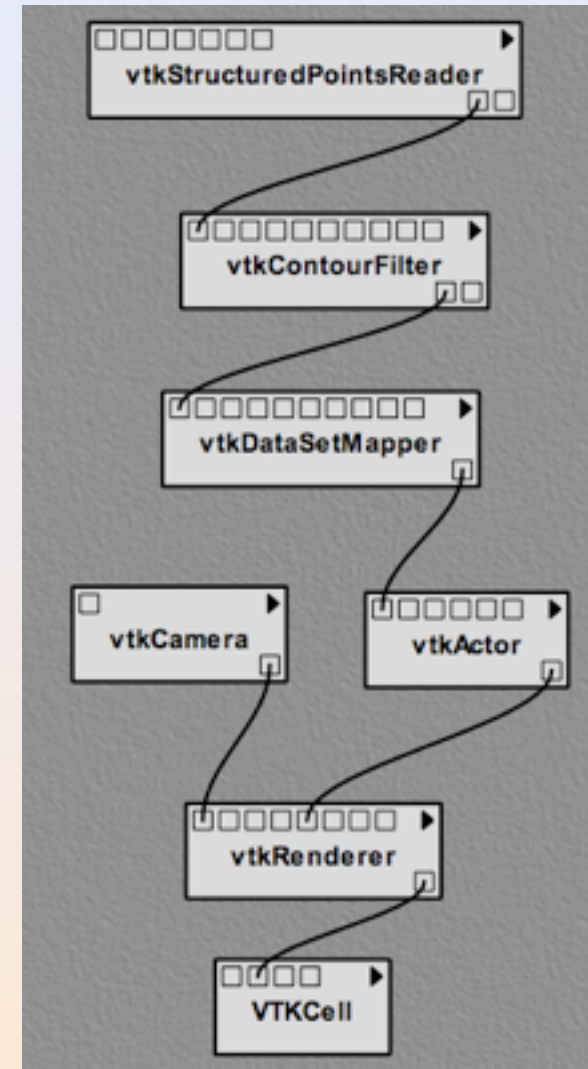
Input:
Head.120.iso

Output:

# Scientific Workflows and Dataflows

- A directed graph describing a computational task
  - Vertices = modules = processing steps + parameters
  - Edges = connections between output and input ports
  - Execution order determined by flow of data from output to input ports
- No state or side effects: Outputs are a *function* of the inputs
- Simple programming model
  - Good match for visual programming interfaces
  - Widely used: adopted by most scientific workflow and visualization systems
  - Easy to optimize and parallelize

I1

vtkStructuredPointsReader
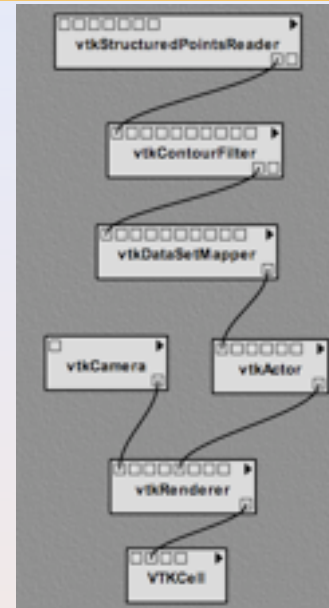
O1

vtkContourFilter

O2

vtkDataSetMapper

O3

# Workflows and Computer Programs

```
1   import vtk
2
3   data = vtk.vtkStructuredPointsReader()
4   data.SetFileName("../examples/data/head.120.vtk")
5
6   contour = vtk.vtkContourFilter()
7   contour.SetInput(0,data.GetOutput())
8   contour.SetValue(0, 67)
9
10  mapper = vtk.vtkPolyDataMapper()
11  mapper.SetInput(contour.GetOutput())
12  mapper.ScalarVisibilityOff()
13
14  actor = vtk.vtkActor()
15  actor.SetMapper(mapper)
16
17  cam = vtk.vtkCamera()
18  cam.SetViewUp(0,0,-1)
19  cam.SetPosition(745,-453,369)
20  cam.SetFocalPoint(135,135,150)
21  cam.ComputeViewPlaneNormal()
22
23  ren = vtk.vtkRenderer()
24  ren.AddActor(actor)
25  ren.SetActiveCamera(cam)
26  ren.ResetCamera()
27
28  renwin = vtk.vtkRenderWindow()
29  renwin.AddRenderer(ren)
30
31  style = vtk.vtkInteractorStyleTrackballCamera()
32  iren = vtk.vtkRenderWindowInteractor()
33  iren.SetRenderWindow(renwin)
34  iren.SetInteractorStyle(style)
35  iren.Initialize()
36  iren.Start()
```

# Workflows and Computer Programs

Program → Workflow

Document → Database



```
<Book>

<Title>The Advanced Html Companion</Title>

<Author> Keith Schengili-Roberts </Author>

<Author> Kim Silk-Copeland</Author>

<Price> 35.96</Price>…

</Book>
```

2. **The Advanced Html Companion**
   by Keith Schengili-Roberts, Kim Silk-Copeland. Paperback (August 1998)
   Our Price:$35.96          Usually ships in 24 hours
   You Save: $8.99 (20%)     Average Customer Review: ★★★★☆

3. **Applied XML Solutions (Sams Professional Publishing)**
   by Benoit Marchal. Paperback (August 29, 2000)
   Our Price:$35.99          Usually ships in 24 hours
   You Save: $9.00 (20%)     Average Customer Review: ★★★☆☆

4. **Applied XML: A Toolkit for Programmers**
   by Alex Ceponkus, Faraz Hoodbhoy. Paperback (July 1, 1999)
   Our Price:$39.99          Usually ships in 24 hours
   You Save: $10.00 (20%)    Average Customer Review: ★★★★☆

*A program is to a workflow what an unstructured document is to a (structured) database.*

# Back to VIS

# ParaView

# VisIt

# VisTrails Project

# Provenance in Art



**Rembrandt van Rijn**
*Self-Portrait, 1659*
*Andrew W. Mellon Collection*
*1937.1.72*

## Provenance

George, 3rd Duke of Montagu and 4th Earl of Cardigan [d. 1790], by 1767;[1] by inheritance to his daughter, Lady Elizabeth, wife of Henry, 3rd Duke of Buccleuch of Montagu House, London; John Charles, 7th Duke of Buccleuch; (P. & D. Colnaghi & Co., New York, 1928); (M. Knoedler & Co., New York); sold January 1929 to Andrew W. Mellon, Pittsburgh and Washington, D.C.; deeded 28 December 1934 to The A.W. Mellon Educational and Charitable Trust, Pittsburgh; gift 1937 to NGA.

[1] This early provenance is established by presence of a mezzotint after the portrait by R. Earlom (1743-1822), dated 1767. See John Charrington, *A Catalogue of the Mezzotints After, or Said to Be After, Rembrandt*, Cambridge, 1923, no. 49.

## Associated Names

- Buccleuch, Henry, 3rd Duke of
- Buccleuch, John Charles, 7th Duke of
- Colnaghi & Co., Ltd., P. & D.
- Knoedler & Company, M.
- Mellon, Andrew W.

# Provenance in Science

- *Provenance is as (or more!) important as the result*
- Not a new issue
- Lab notebooks have been used for a long time
- What is new?
  - Large volumes of data
  - Complex analyses
- Writing notes is no longer an option…
- Computational provenance



*DNA recombination By Lederberg*

# Provenance in Science

- *Provenance is as (or more!) important as the result*
- Not a new issue
- Lab notebooks have

been used for a long time

- What is new?
  - Large volumes of data
  - Complex analyses
- Writing notes is no longer an option…
- Computational provenance

*When*



*DNA recombination By Lederberg*

# Provenance in Science

- *Provenance is as (or more!) important as the result*
- Not a new issue
- Lab notebooks have been used for a long time
- What is new?
  - Large volumes of data
  - Complex analyses
- Writing notes is no longer an option…
- Computational provenance

*When*

*Observed data*

*DNA recombination By Lederberg*

# Provenance in Science

- *Provenance is as (or more!) important as the result*
- Not a new issue
- Lab notebooks have been used for a long time
- What is new?
  - Large volumes of data
  - Complex analyses
- Writing notes is no longer an option…
- Computational provenance

*When*

*Annotation*

*Observed data*

*DNA recombination By Lederberg*

# Exploration and Workflows

- Workflows have been traditionally used to automate repetitive tasks
- In exploratory tasks, *change is the norm*!
  - Data analysis and exploration are iterative processes



Figure modified from J. van Wijk, IEEE Vis 2005

# Exploration and Creativity Support

- Reflective reasoning is key in the exploratory processes
- "*Reflective reasoning requires the ability to store temporary results, to make inferences from stored knowledge, and to follow chains of reasoning backward and forward, sometimes backtracking when a promising line of thought proves to be unfruitful. …the process is slow and laborious*"

  Donald A. Norman

- Need external aids—tools to facilitate this process
  - Creativity support tools [Shneiderman, CACM 2002]
- Need aid from people—collaboration

# The need for provenance

*What's the difference?*

anon4877_base_20060331.jpg          anon4877_lesion_20060401.jpg



How were these images created?

Are they really from the same patient?

# The need for provenance

anon487**6**_base_20060331.jpg          anon4877_lesion_20060401.jpg



How were these images created?

Are they really from the same patient?

# Data Exploration and Workflows

# Data Exploration and Workflows

raw data:CT scan

workflow

# Data Exploration and Workflows

workflow

raw data:CT scan

# Data Exploration and Workflows

raw data:CT scan

workflow

# Data Exploration and Workflows



raw data:CT scan

workflow

Files (workflow specifications)
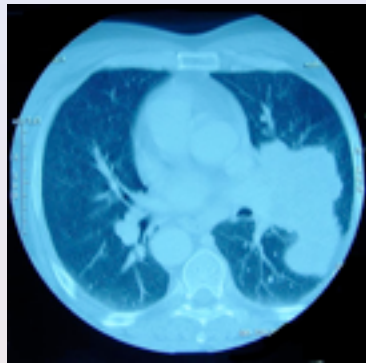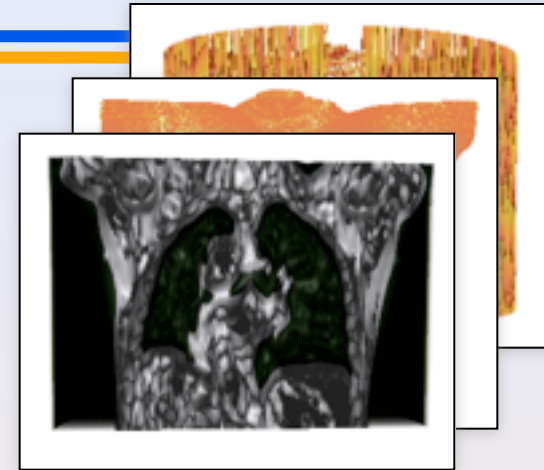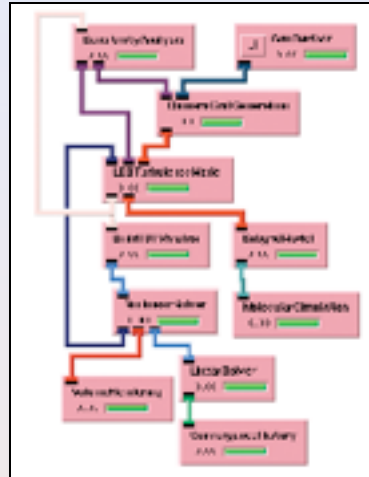
anon4877_voxel_scale_1_zspace_20060331.srn

Notes

Initial visualization with z-scaling corrected

# Data Exploration and Workflows



raw data:CT scan

workflow

Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

Notes

Initial visualization with z-scaling corrected

# Data Exploration and Workflows

raw data:CT scan

workflow



Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

Notes

Initial visualization with z-scaling corrected

# Data Exploration and Workflows



raw data:CT scan

workflow

Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

Notes

Initial visualization wi

Added texture and shading

# Data Exploration and Workflows
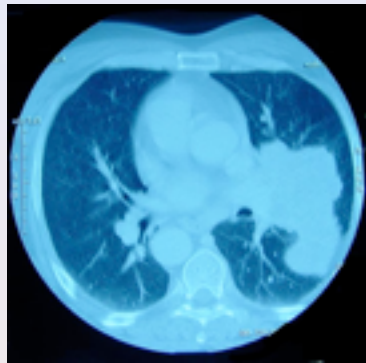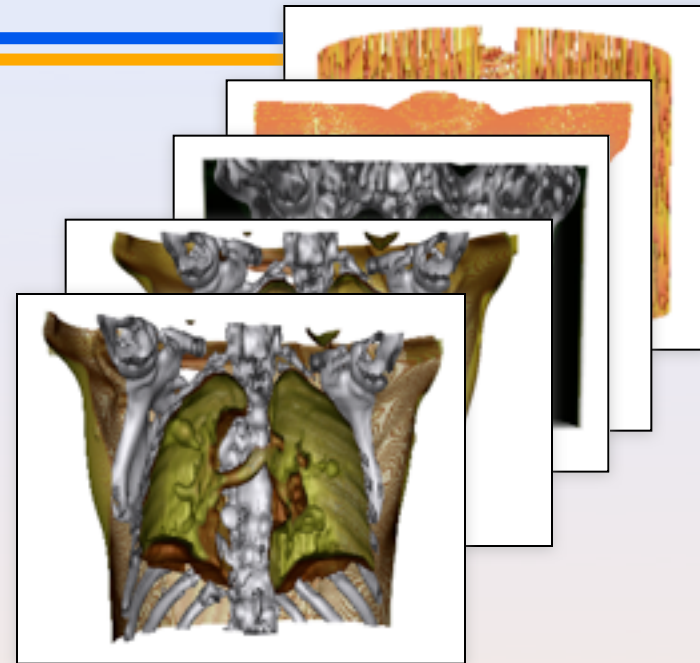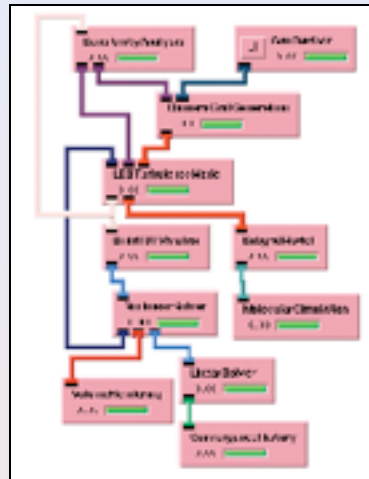
raw data:CT scan

workflow



Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

Notes

Initial
visualization
wi... Added texture
and shading

# Data Exploration and Workflows

**workflow**

raw data:CT scan



## Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

anon4877_textureshading_plane0_20060331.srn

## Notes

Initial visualization wi...

Added texture an...

Added plane to visualize internal structure

# Data Exploration and Workflows

**workflow**

**raw data:CT scan**



## Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

anon4877_textureshading_plane0_20060331.srn

## Notes

Initial visualization wi...

Added texture an...

Added plane to visualize internal structure

# Data Exploration and Workflows

workflow

raw data:CT scan



Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

anon4877_textureshading_plane0_20060331.srn

anon4877_goodxferfunction_20060331.srn

Notes

Initial visualization wi...

Added texture a...

Added plane to v...

Found good transfer function

# Data Exploration and Workflows

workflow

raw data:CT scan



Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

anon4877_textureshading_plane0_20060331.srn

anon4877_goodxferfunction_20060331.srn

Notes

Initial visualization wi...

Added texture a...

Added plane to v... i... s...

Found good transfer function

# Data Exploration and Workflows

**workflow**

raw data:CT scan



## Files (workflow specifications)

anon4877_voxel_scale_1_zspace_20060331.srn

anon4877_textureshading_20060331.srn

anon4877_textureshading_plane0_20060331.srn

anon4877_goodxferfunction_20060331.srn

anon4877_lesion_20060331.srn

## Notes

Initial visualization wi...

Added texture a...

Added plane to v...

Found good ...

Identified lesion tissue

# VisTrails: Managing Exploration

- Comprehensive *provenance infrastructure* for computational tasks
  - Data + workflow provenance
  - *Treat workflow as a 1st-class data product*
- Support for *exploratory* tasks such as visualization and data mining
  - Task specification iteratively refined as users generate and test hypotheses
- VisTrails manages the data, metadata and the exploration process, scientists can focus on *science!*
- Not a replacement for visualization or scientific workflow systems: infrastructure that can be combined with and enhance these systems
- Focus on usability—build tools for scientists

http://www.vistrails.org

# Keeping Exploration Trails



**Trail**

# Keeping Exploration Trails



**Trail**

Workflows

Data Products

# Keeping Exploration Trails

**Trail**

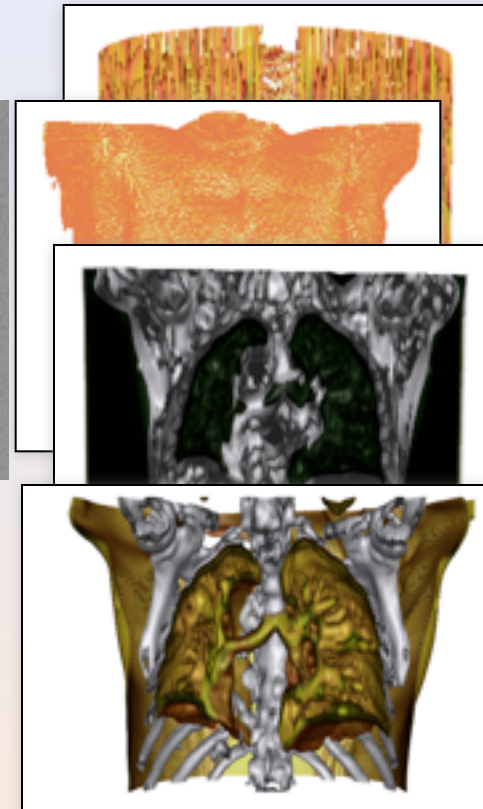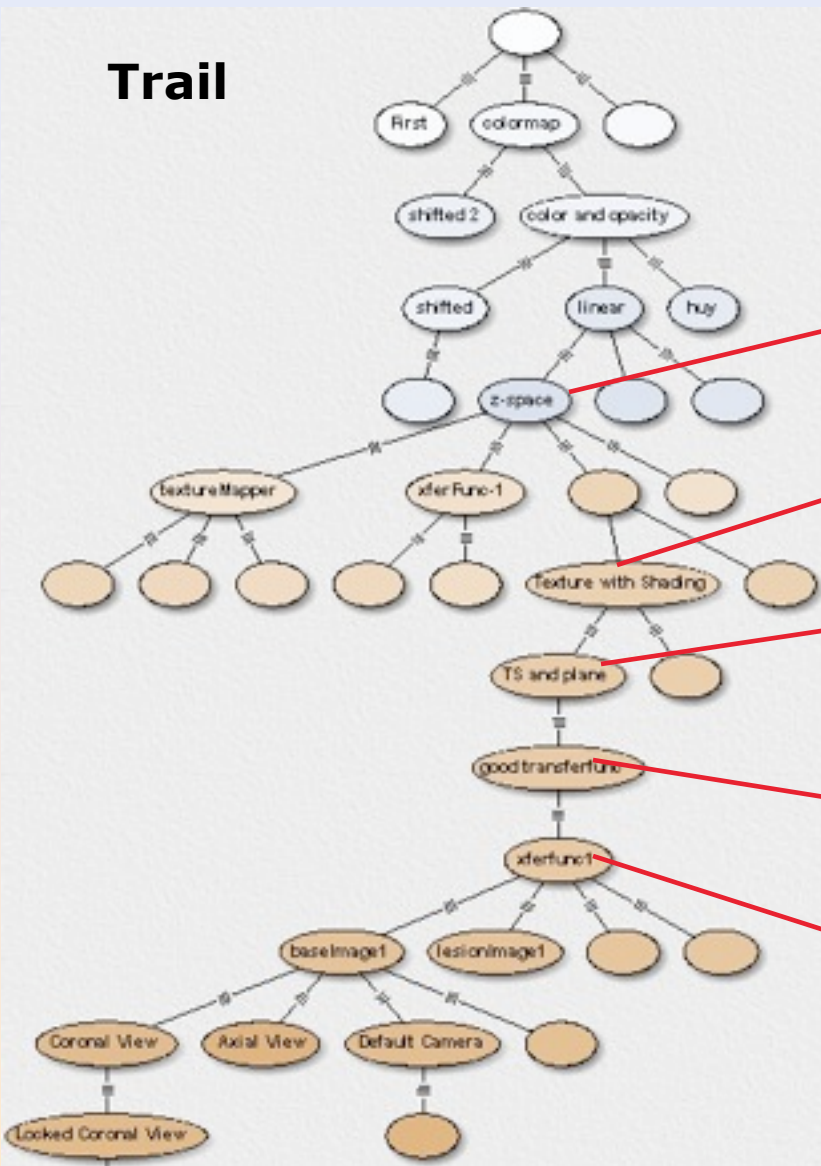

Workflows

Data Products
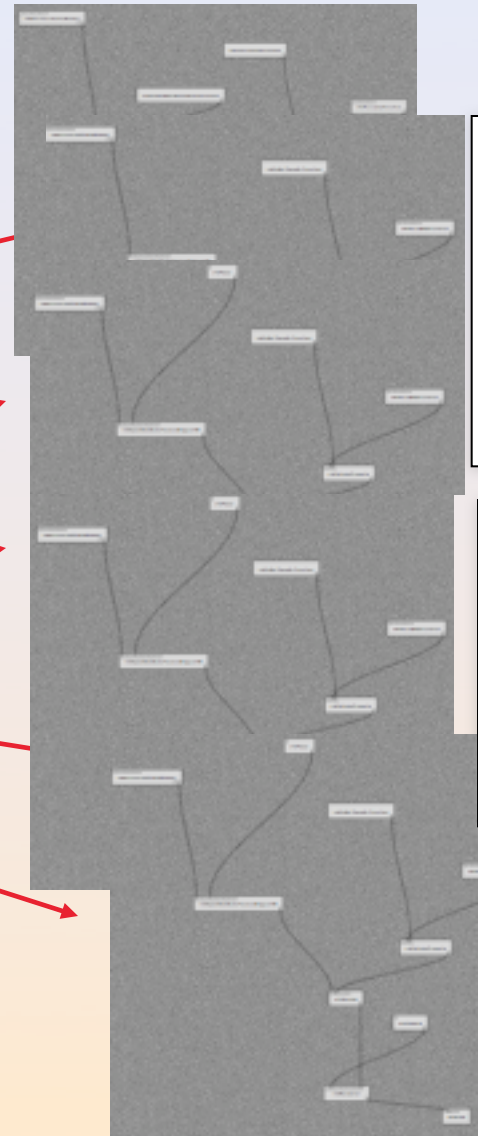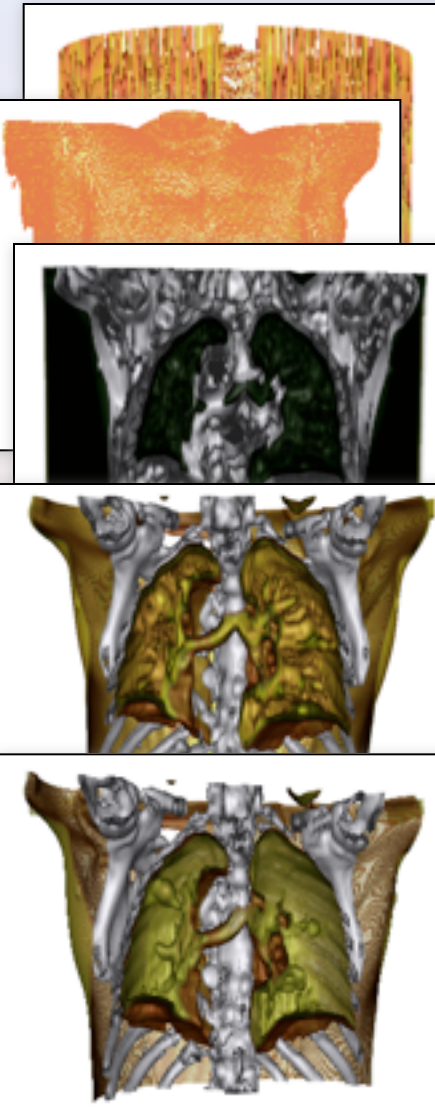
# Keeping Exploration Trails



**Trail**

Workflows

Data Products

# Keeping Exploration Trails



**Trail**

Workflows

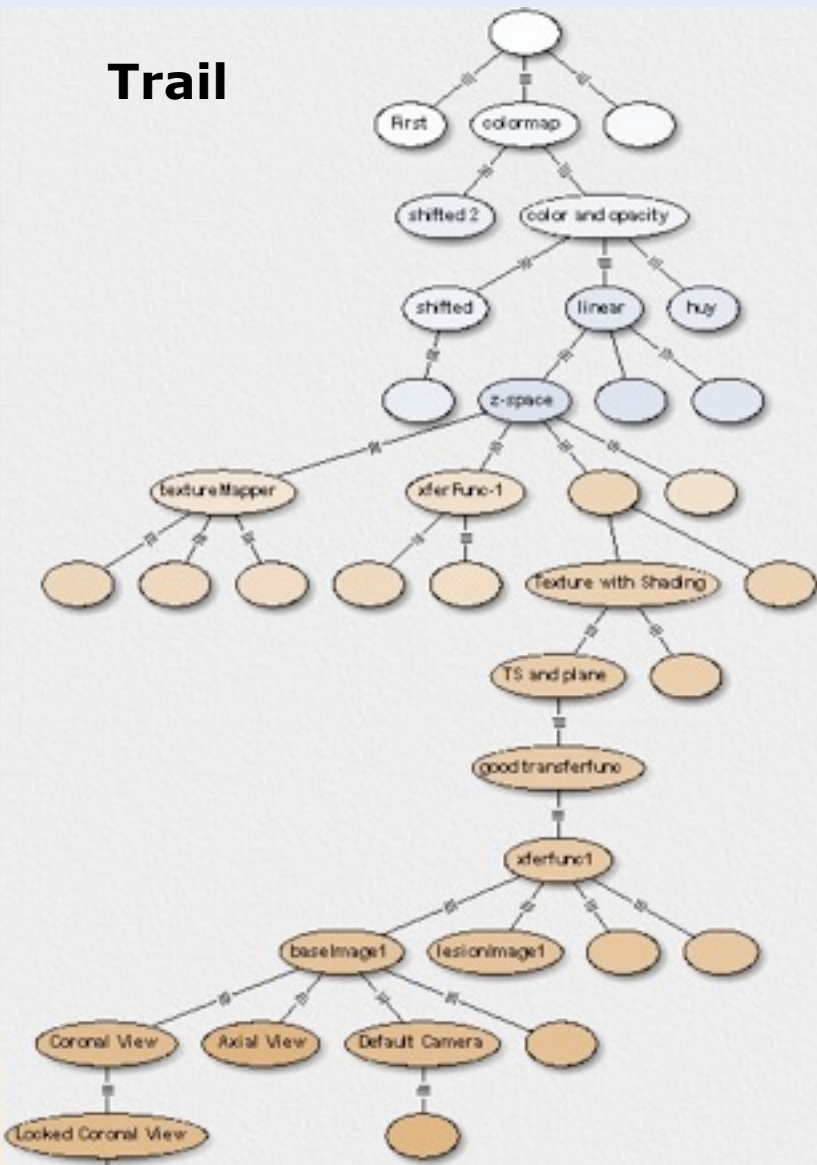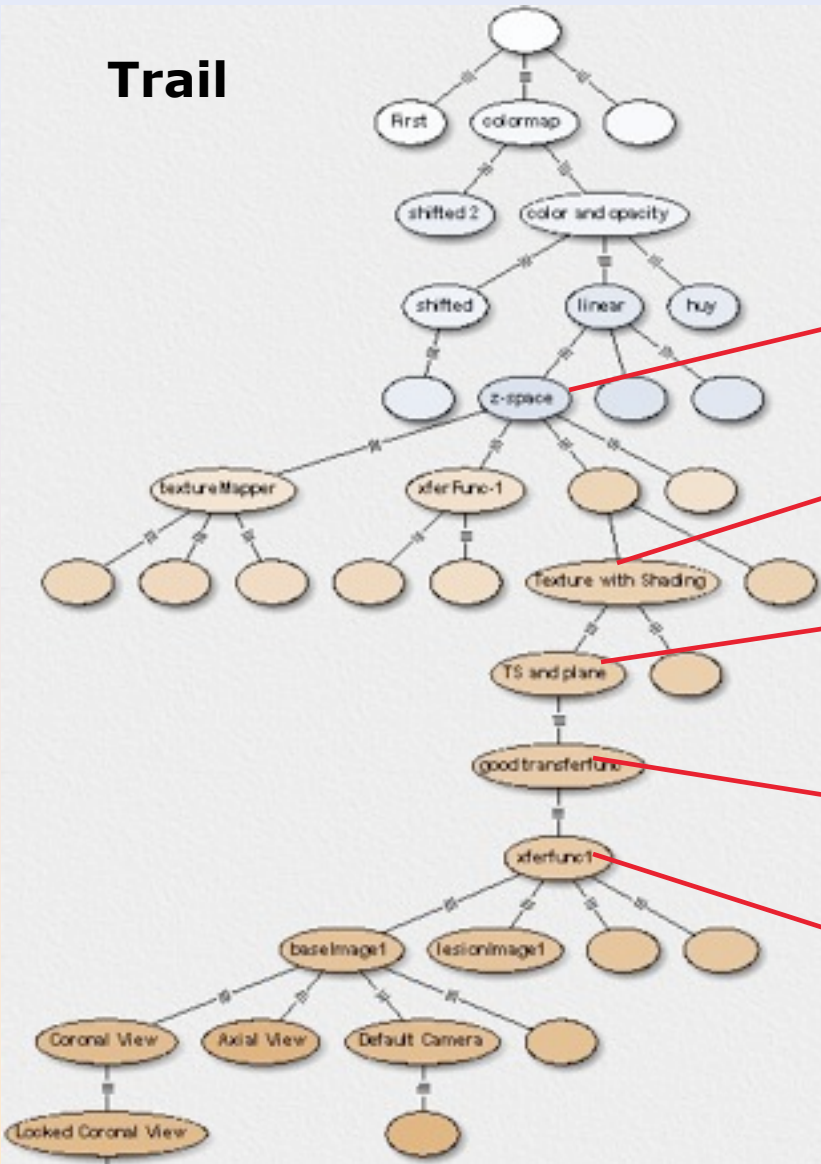Data Products

# Keeping Exploration Trails

**Trail**



Workflows

Data Products

# Keeping Exploration Trails



**Trail**

Workflows

Data Products

Software for Exploratory Visualization

# Keeping Exploration Trails



**Trail**

# Keeping Exploration Trails



**Trail**

Notes

Initial visualization with z-scaling corrected

Added texture and shading

Added plane to visualize internal structure

Found good transfer function

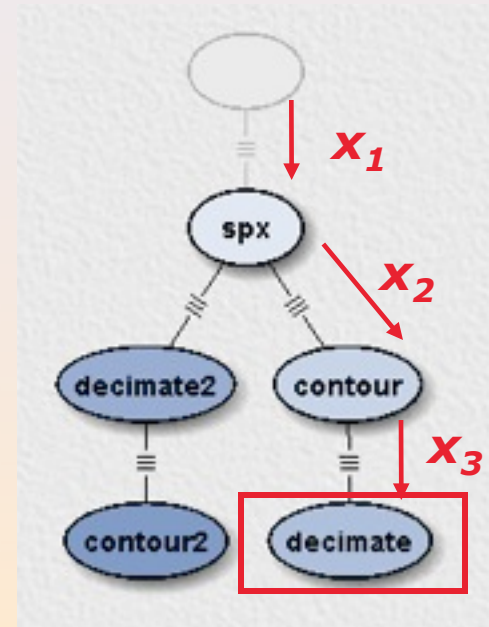Identified lesion tissue

# Keeping Exploration Trails

# Change-Based Provenance



- Records actions
- Provenance = changes to computational tasks
  - Add a module, add a connection, change a parameter value
- Extensible *change* algebra

addModule

deleteConnection

addConnection

addConnection

setParameter

# Change-Based Provenance

- ☒ Records actions
- ☒ Provenance = changes to computational tasks
  - Add a module, add a connection, change a parameter value
- ☒ Extensible *change* algebra
- ☒ A *vistrail* node $v_t$ corresponds to the workflow that is constructed by the sequence of actions from the root to $v_t$

$$v_t = x_n \circ x_{n-1} \circ \dots \circ x_1 \circ \varnothing$$

[Freire et al, IPAW 2006]

*vistrail*

# Change-based provenance

- Records changes to workflows
- Workflow evolution is captured in a vistrail, a rooted tree where

–nodes correspond to workflow versions

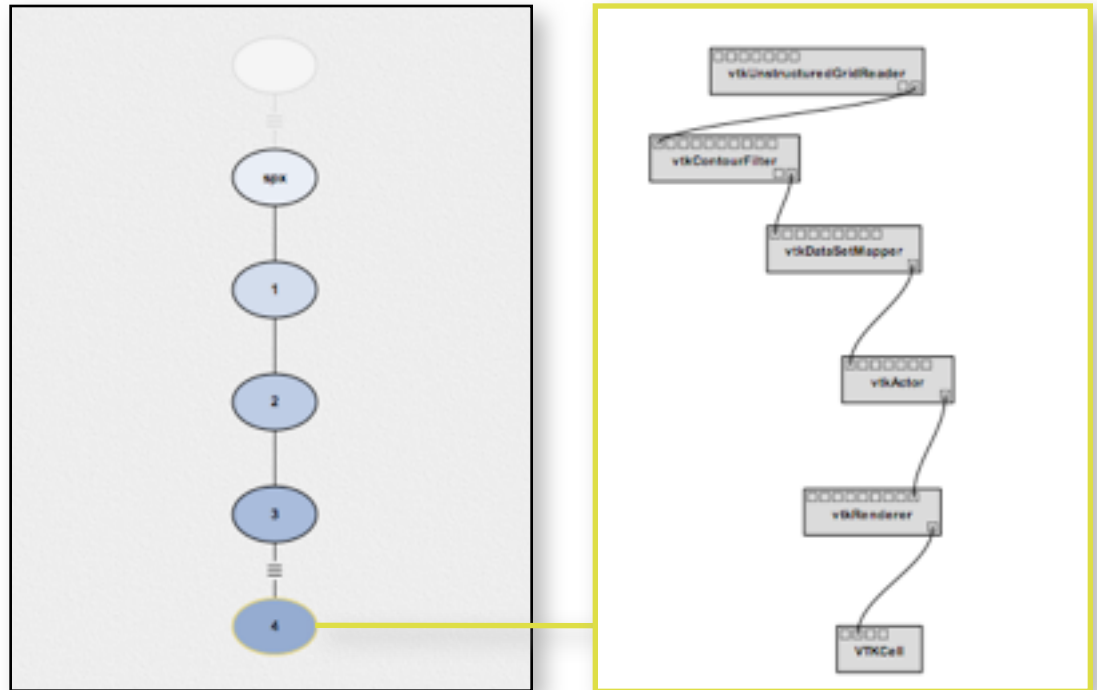–edges correspond to actions that transform the parent into the child workflow

# Change-based provenance

# Change-based provenance
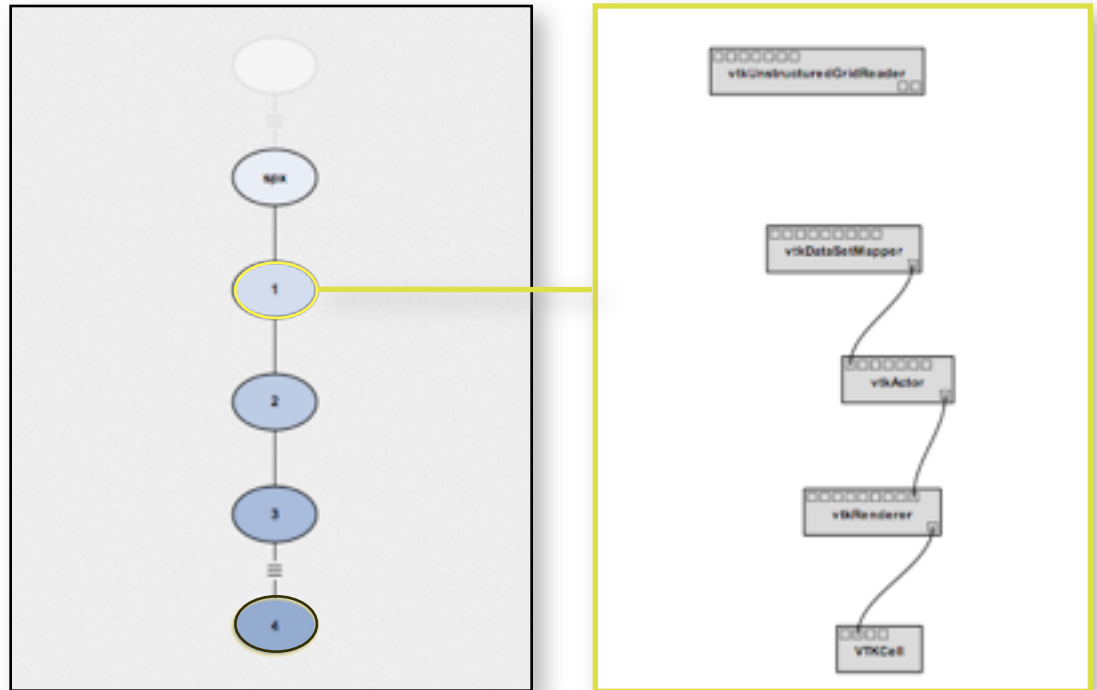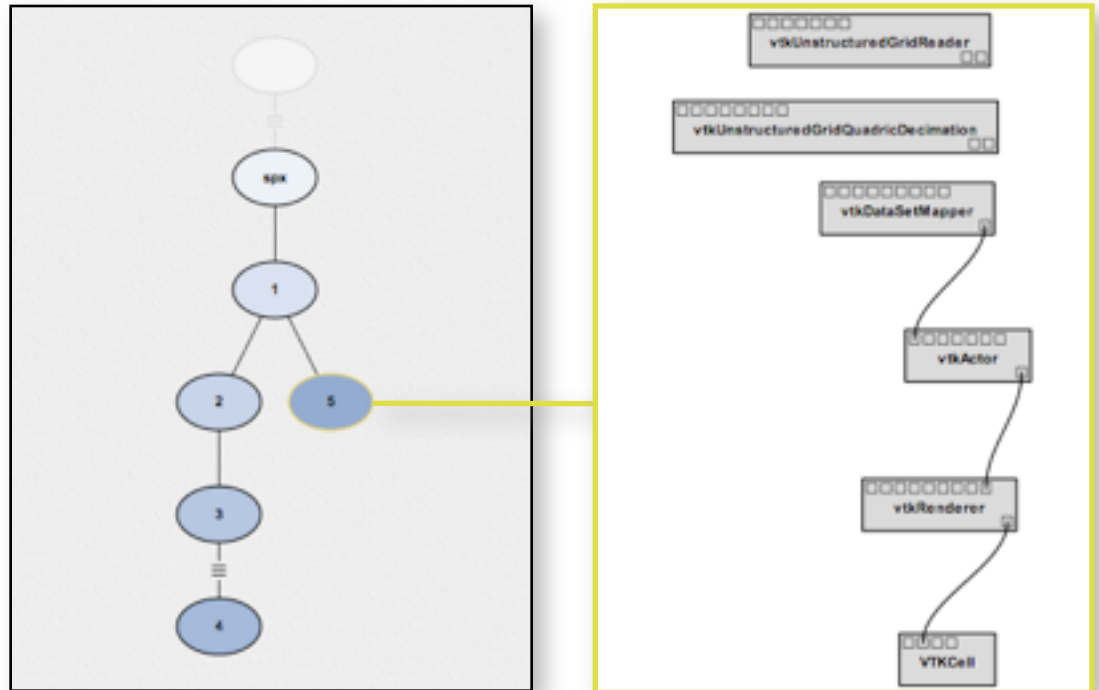
# Change-based provenance

# Change-based provenance

# Change-based provenance

# Change-based provenance

# Change-Based Provenance

- General: Works with any system that has undo/redo!
- Concise representation
- Uniformly captures data and workflow provenance
  - Data provenance: where does a specific data product come from?
  - Workflow evolution: how has workflow structure changed over time?
- Detailed information about the exploration process
  - Results can be reproduced
- Provenance beyond reproducibility:
  - Scientists can return to any point in the exploration space
  - Enables scalable exploration of the parameter space
  (and compare results using a spreadsheet)
  - Support for collaboration
  - Understand problem-solving strategies—knowlegde re-use

# What's the difference?
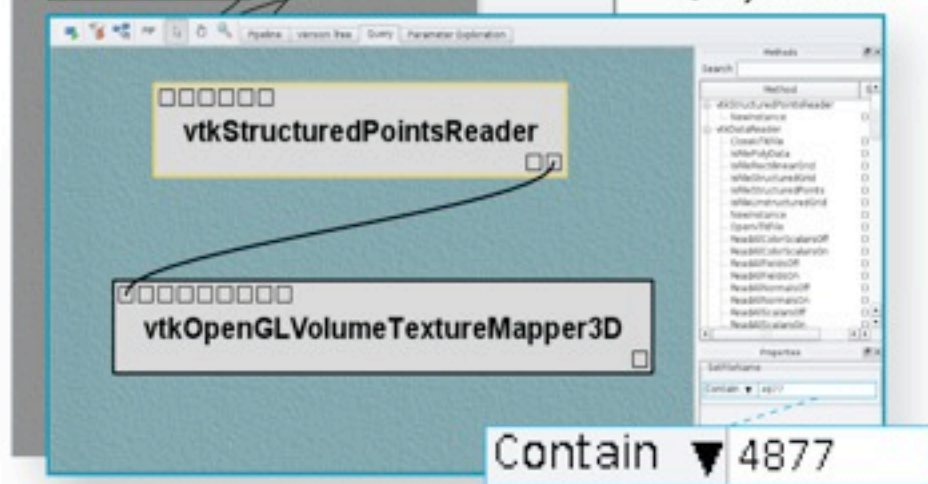


baseImage1

lesionImage1

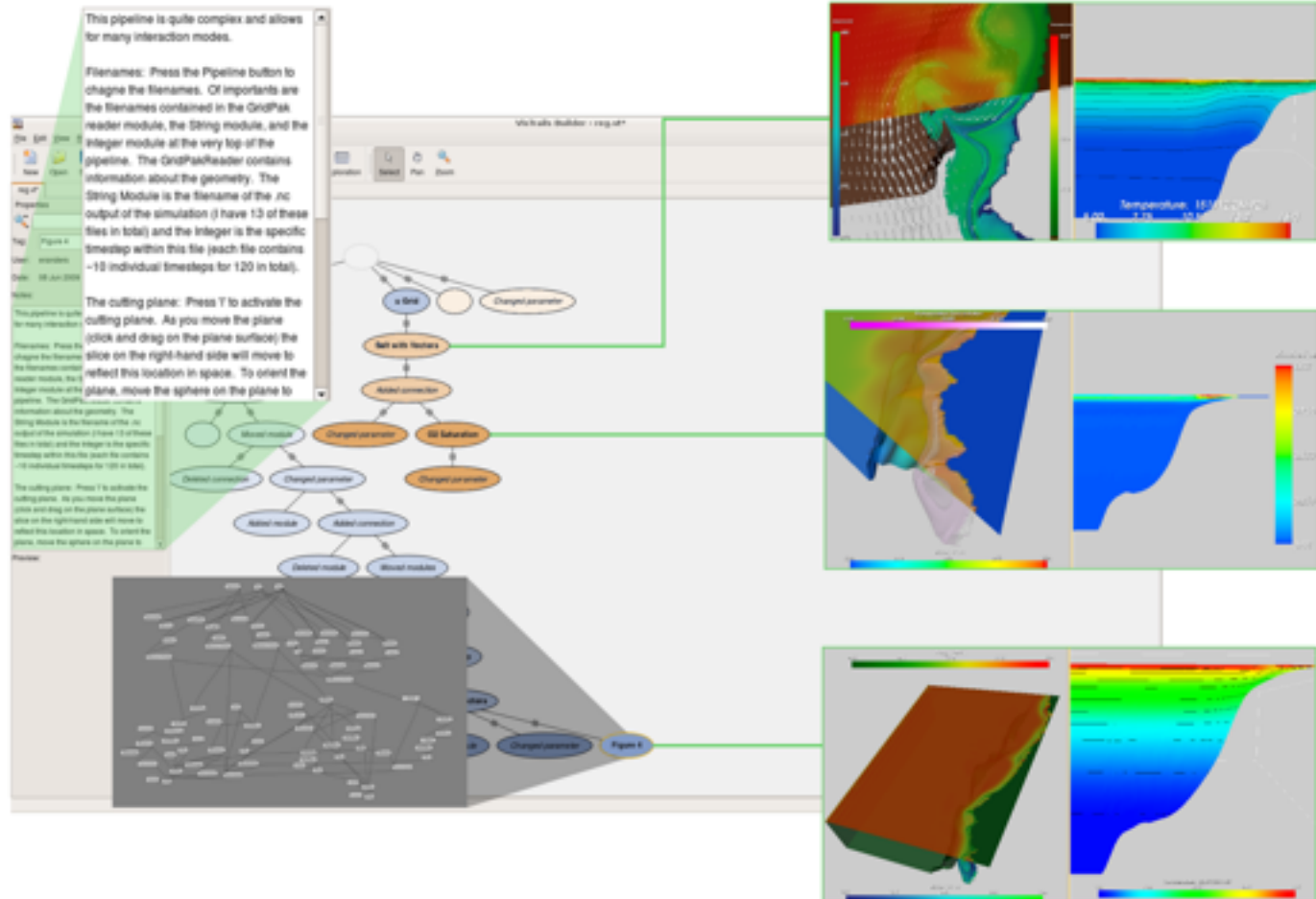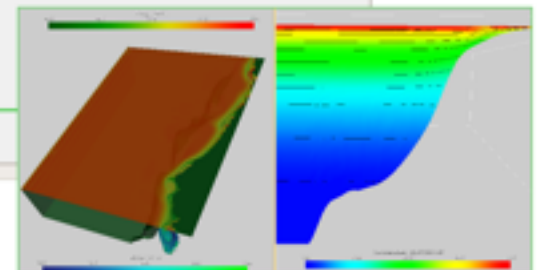# Differences in Specification

# Query by Example

# VisTrails: Provenance of Exploration

Reproducibility
and
Validation

# VisTrails: Provenance of Exploration
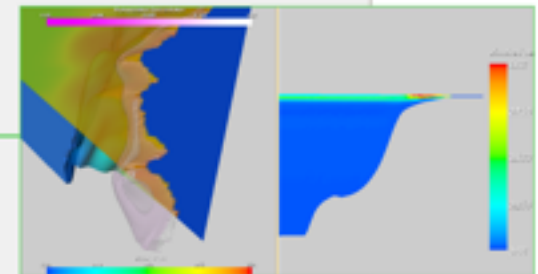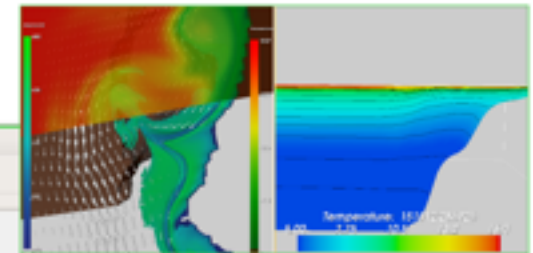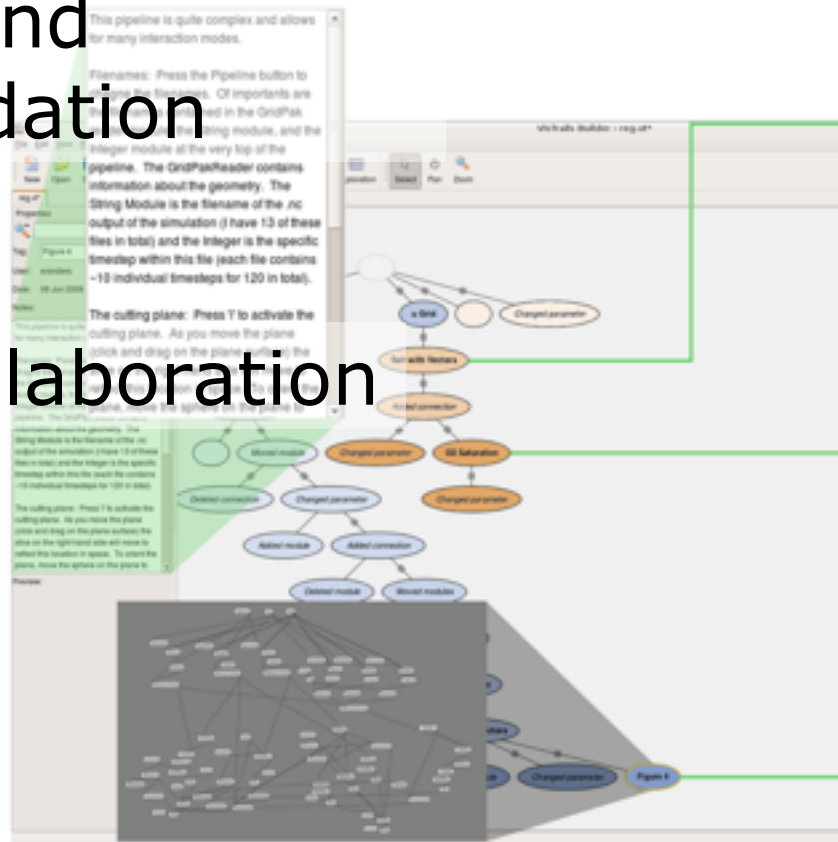
Reproducibility and Validation

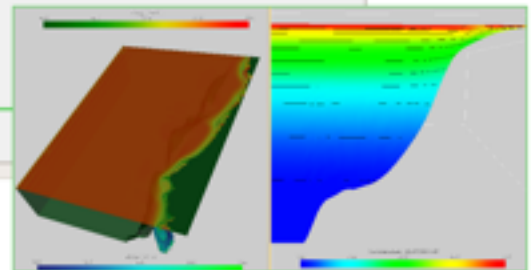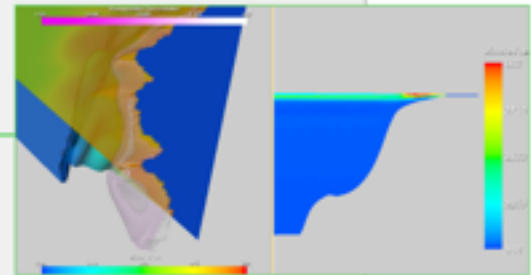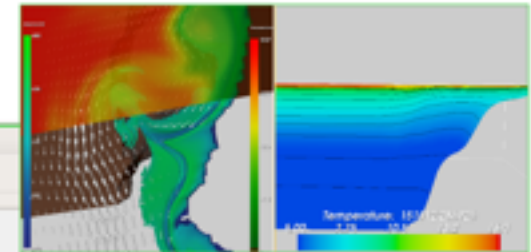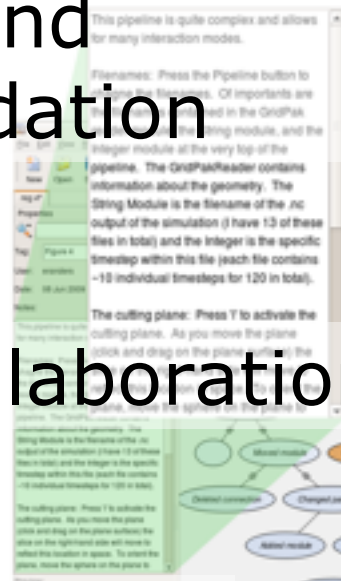Collaboration

# VisTrails: Provenance of Exploration
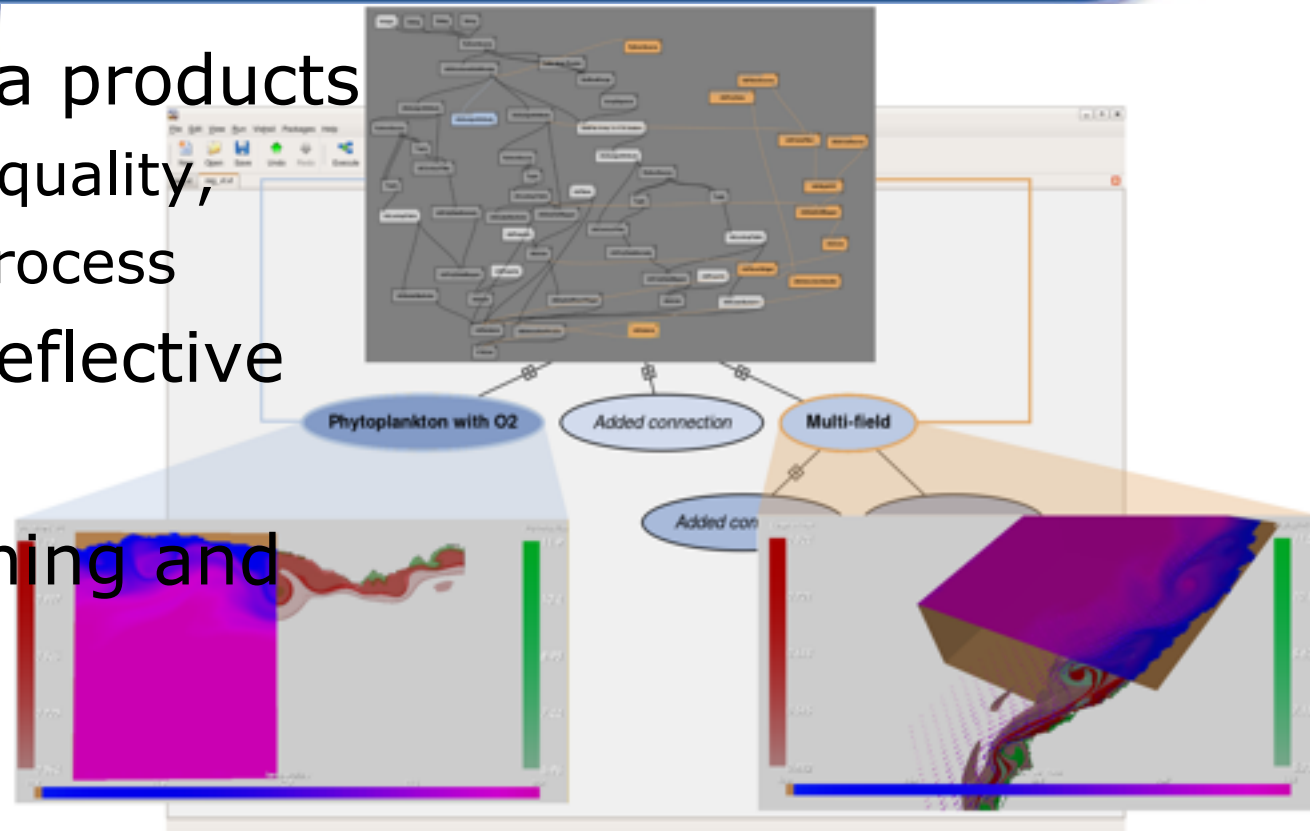
Reproducibility and Validation

Collaboration

Knowledge sharing: data + processes

# Benefits of Provenance

- Compare data products
  - Assess data quality, understand process
- Support for reflective reasoning
- Improve training and teaching

[Freire et al., IPAW 2006]

# Benefits of Provenance

- Compare data products
  - Assess data quality, understand process
- Support for reflective reasoning
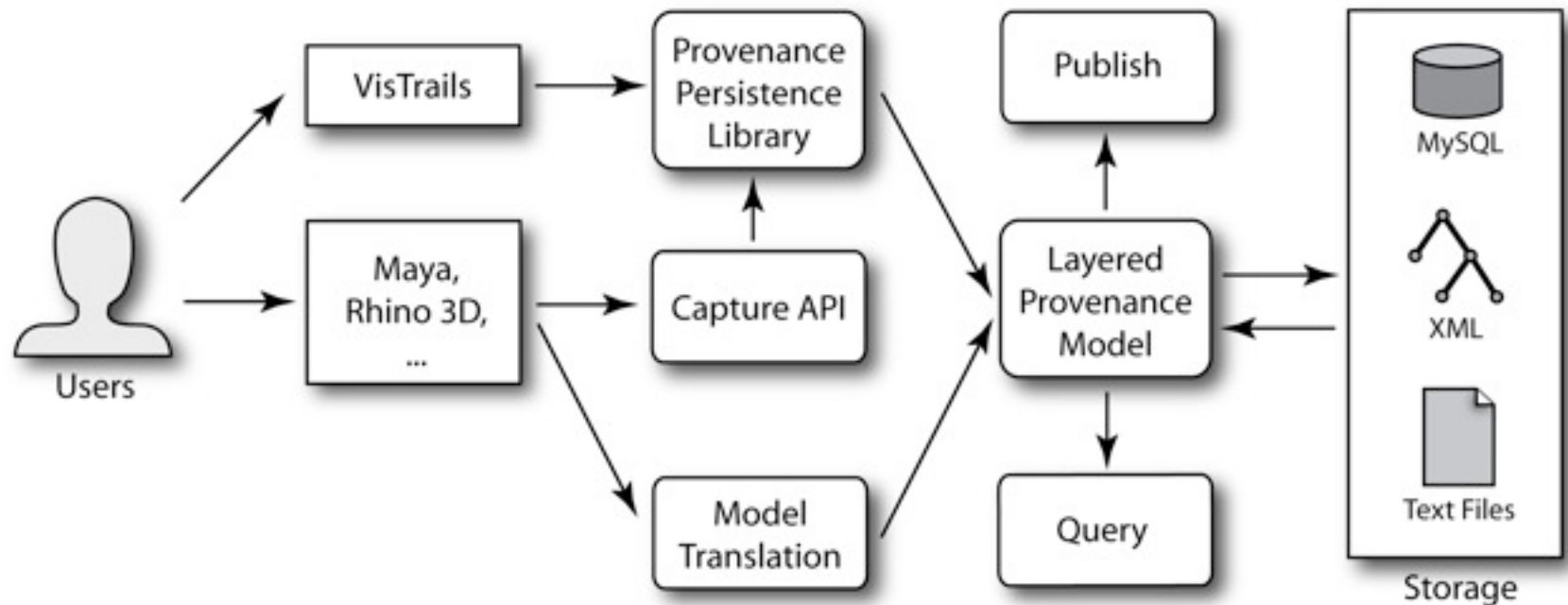- Improve training and teaching



*"Reflective reasoning requires the ability to store temporary results, to make inferences from stored knowledge, and to follow chains of reasoning backward and forward, sometimes backtracking when a promising line of thought proves to be unfruitful. …the process is slow and laborious"*
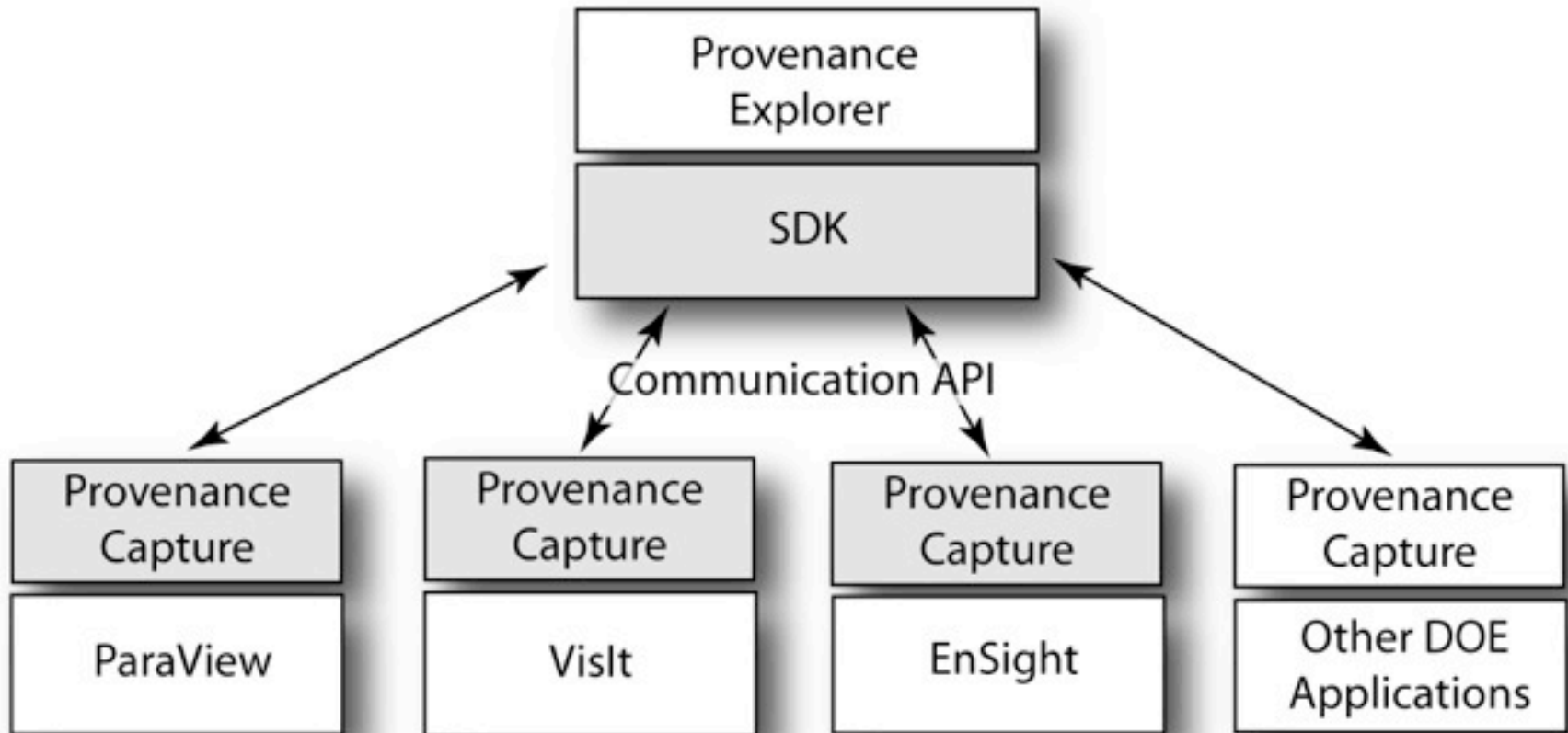Things that make us smart, Donald A. Norman
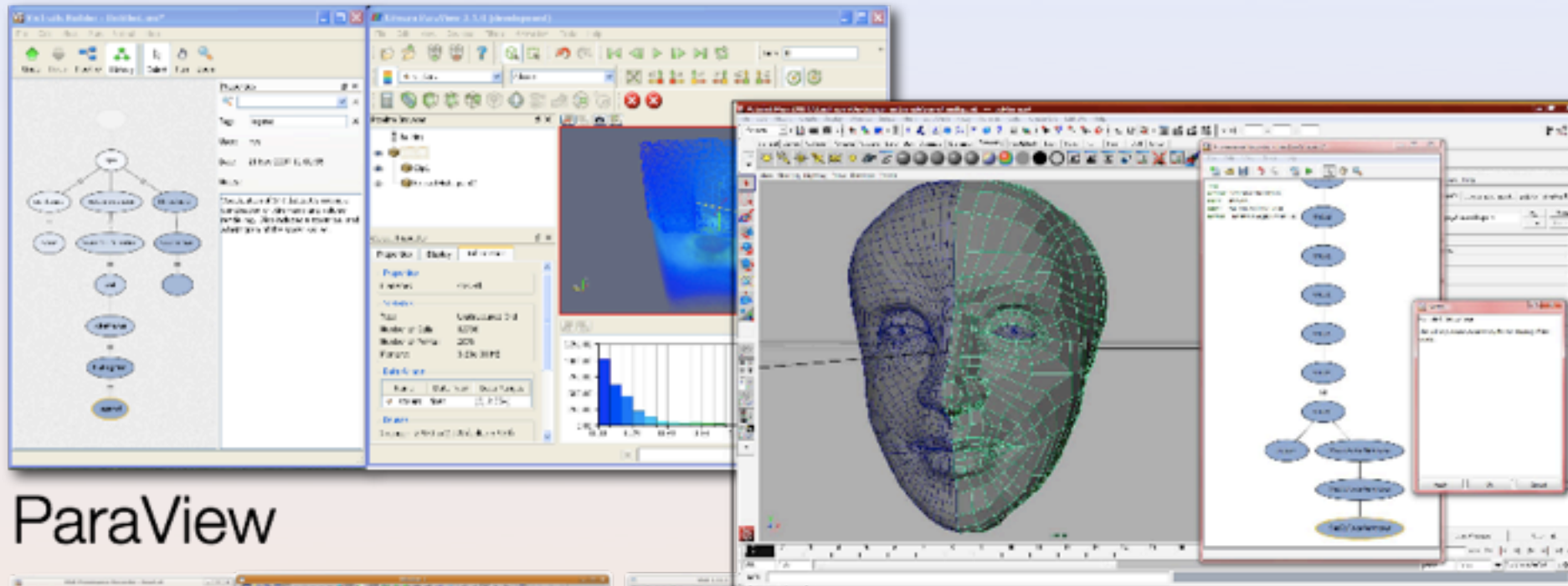
[Freire et al., IPAW 2006]

# Benefits of Provenance

- Compare data products
  - Assess data quality, understand process
- Support for reflective reasoning
- Improve training and teaching

Phytoplankton with O2  Added connection  Multi-field

Added con

*"Reflective reasoning requires the ability to store temporary results, to make inferences from stored knowledge, and to follow chains of reasoning backward and forward, sometimes backtracking when a promising line of thought proves to be unfruitful. ...the process is slow and laborious"*
Things that make us smart, Donald A. Norman

[Freire et al., IPAW 2006]

# Provenance API



**www.vistrails.org**

# Provenance "Plug-ins"
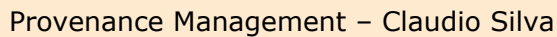
# Provenance Enabling Tools



ParaView

Maya

[Callahan et al., IPAW 2008]

VisIt

# More plugins...

ImageSeg3D
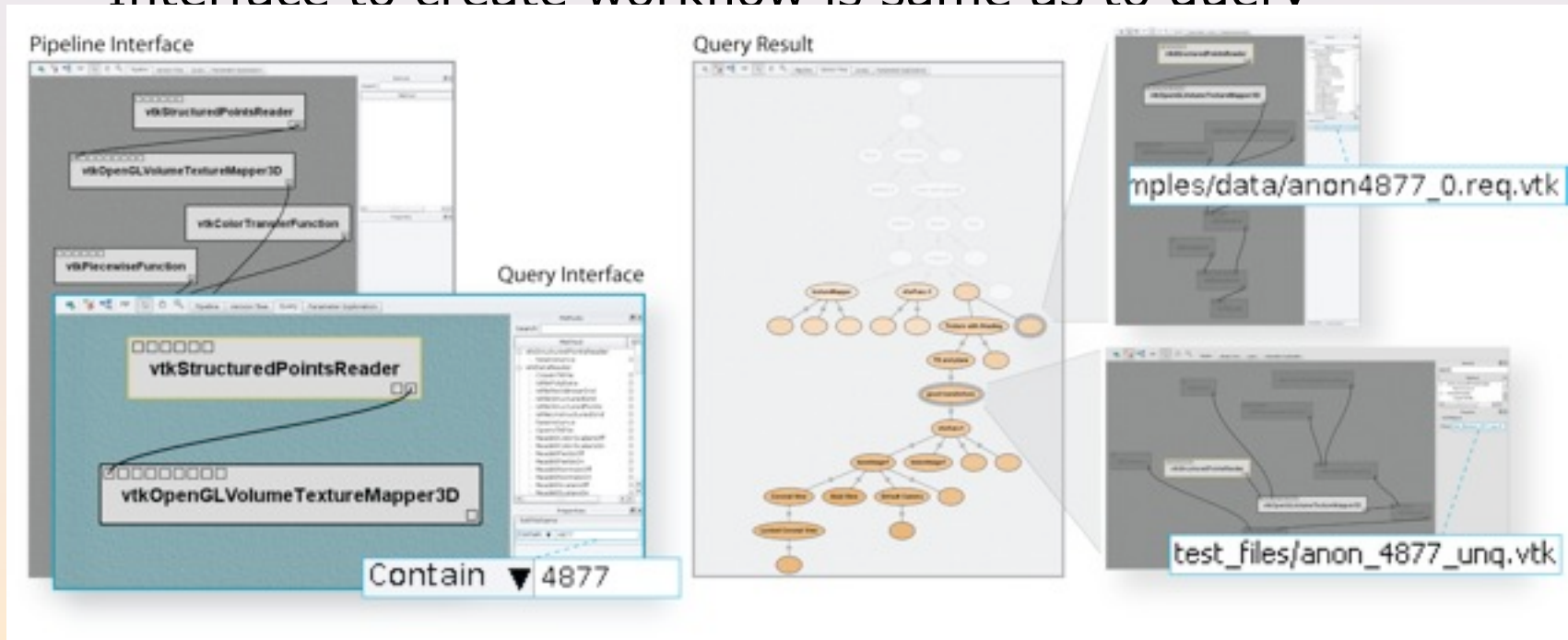
# Sample of Ongoing Work

# Querying Workflows by Example
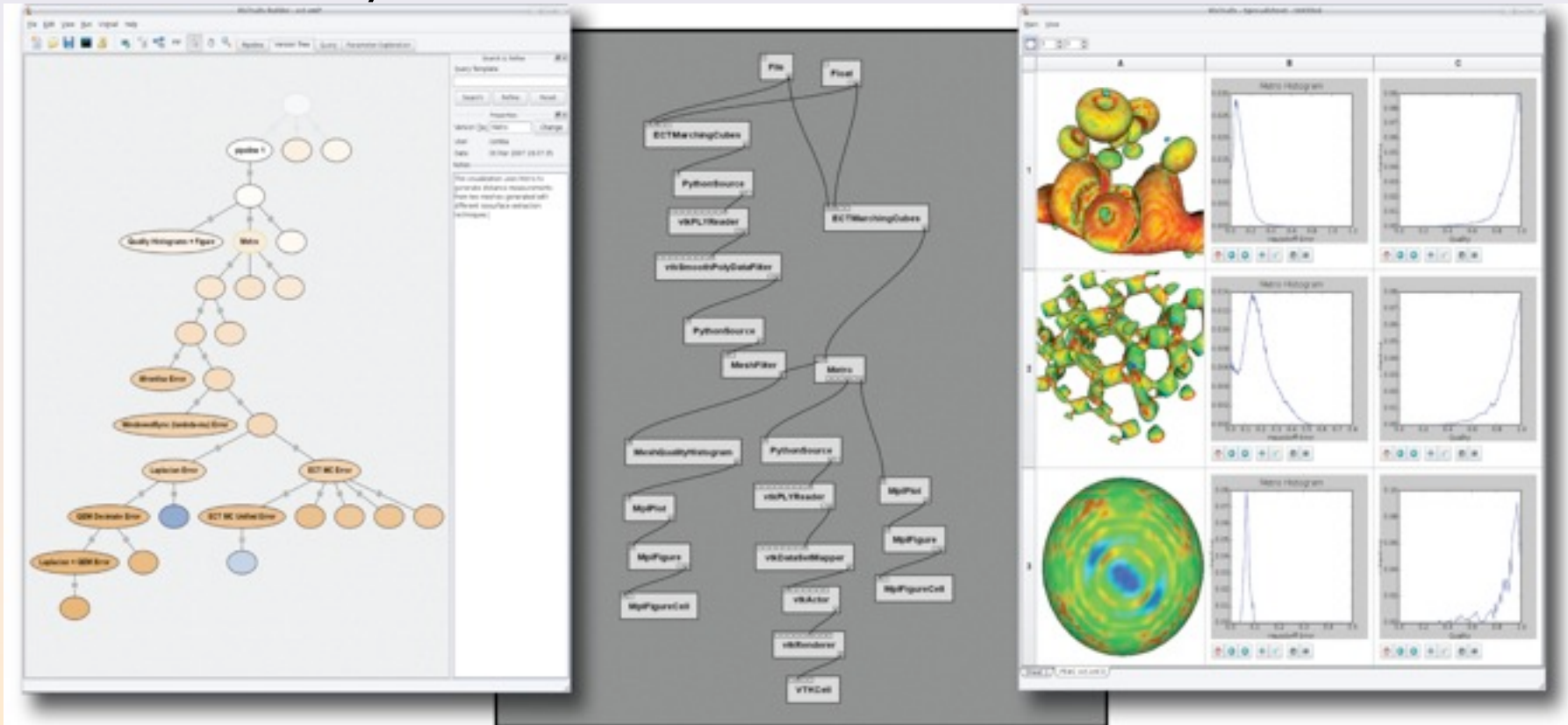
- Workflows are graphs: hard to specify queries using text!
- Querying workflows by example [Scheidegger et al., TVCG 2007; Beeri et al., VLDB 2006; Beeri et al. VLDB 2007]
  - WYSIWYQ -- What You See Is What You Query
  - Interface to create workflow is same as to query

# Creating Workflows

- Complex workflows are hard to create
  - Programming expertise
  - Domain knowledge
  - Familiarity with different tools

# Creating Workflows

☒ Complex workflows are hard to create

– Programming expertise

– Domain knowledge

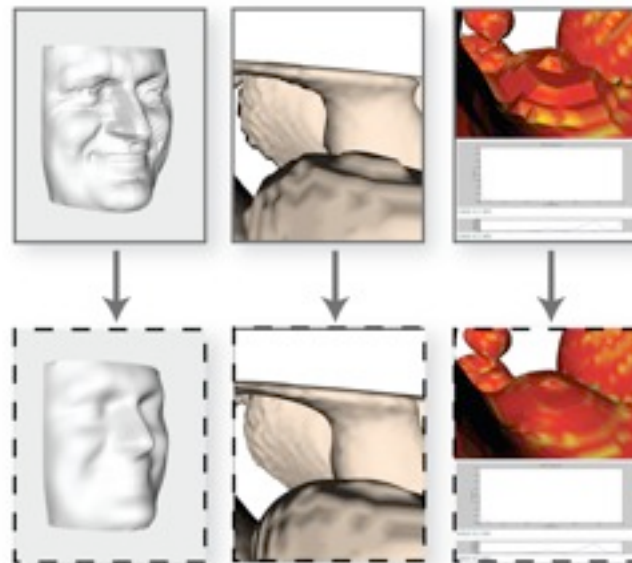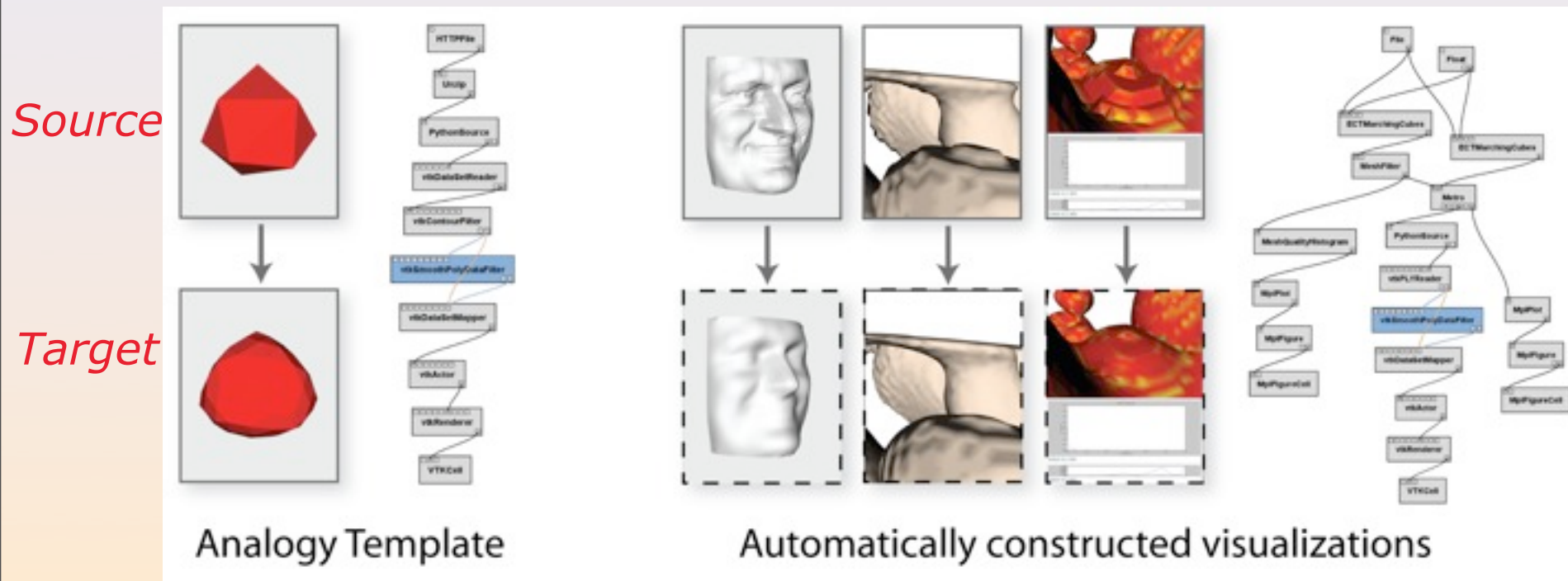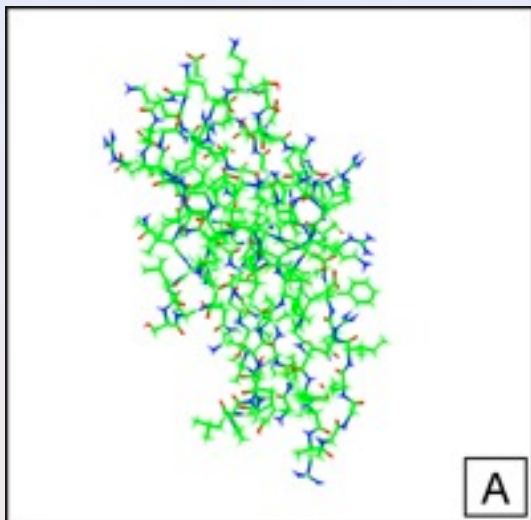*Steep learning curve*

– Familiarity with different tools

# Creating Workflows by Analogy

- Use the wisdom of the crowds
  - Some workflow refinements are common, e.g., change the rendering technique, publish image on the Web
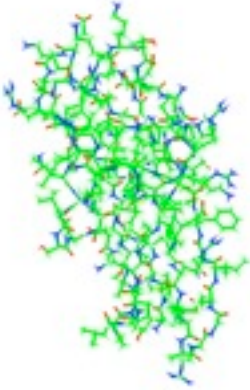- Apply refinements by analogy, automatically [Scheidegger et al, IEEE TVCG 2007]



Analogy Template                Automatically constructed visualizations
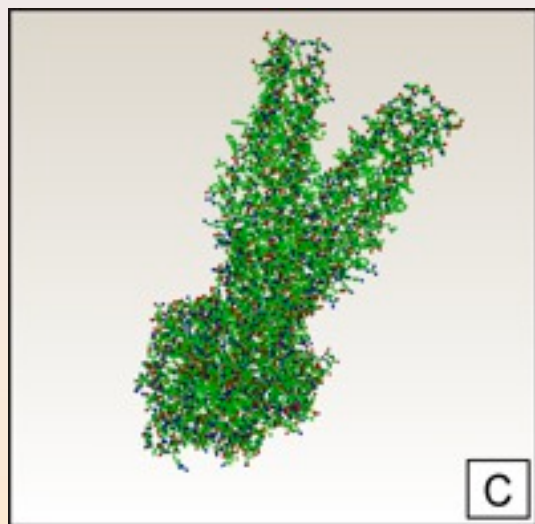
# Creating Workflows by Analogy

- Use the wisdom of the crowds
  - Some workflow refinements are common, e.g., change the rendering technique, publish image on the Web
- Apply refinements by analogy, automatically [Scheidegger et al, IEEE TVCG 2007]

*Source*



Analogy Template                    Automatically constructed visualizations

# Creating Workflows by Analogy

- Use the wisdom of the crowds
  - Some workflow refinements are common, e.g., change the rendering technique, publish image on the Web
- Apply refinements by analogy, automatically [Scheidegger et al, IEEE TVCG 2007]



Source

Target

Analogy Template

Automatically constructed visualizations

# Creating Workflows by Analogy



is to



as



is to

?

# Creating Workflows by Analogy



A **is to** B **as** C **is to** D

# Creating Workflows by Analogy
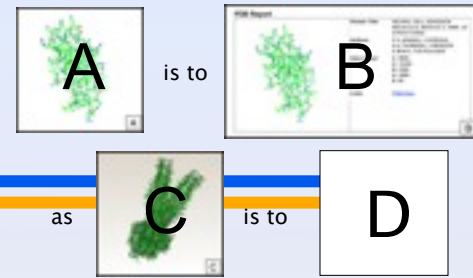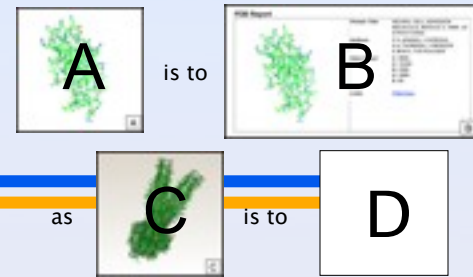


1. Compute difference: Δ(A,B)
   – Just like a patch!
   – But…

   D = Δ(A,B) ∘ C may not be a valid workflow

# Creating Workflows by Analogy



A is to B as C is to D

1. Compute difference: Δ(A,B)
   – Just like a patch!
   – But…

   D = Δ(A,B) ∘ C may not be a valid workflow



$\Delta = DIFF(A, B)$

# Creating Workflows by Analogy

1. Compute difference: Δ(A,B)
   - Just like a patch!
   - But...

   D = Δ(A,B) ◦ C may not be a valid workflow

# Creating Workflows by Analogy

1. Compute difference: Δ(A,B)
   - Just like a patch!
   - But…

   D = Δ(A,B) ◦ C may not be a valid workflow

u. Find correspondences between A and C: map(A,C)
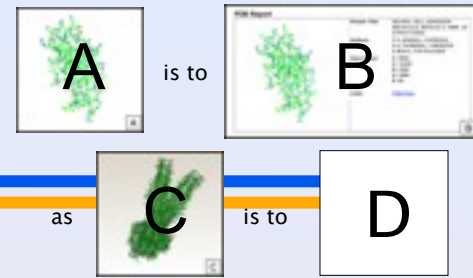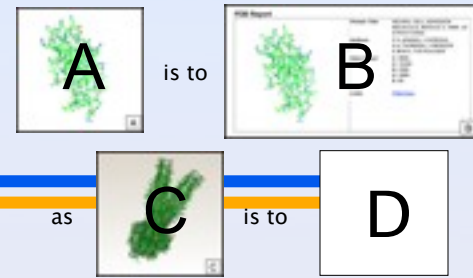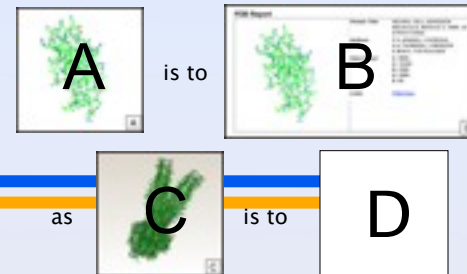   - Diffuse similarity scores across the product graph AxC using Eigenvalue decompositions

# Creating Workflows by Analogy



A is to B as C is to D

1. Compute difference: Δ(A,B)
   – Just like a patch!
   – But…

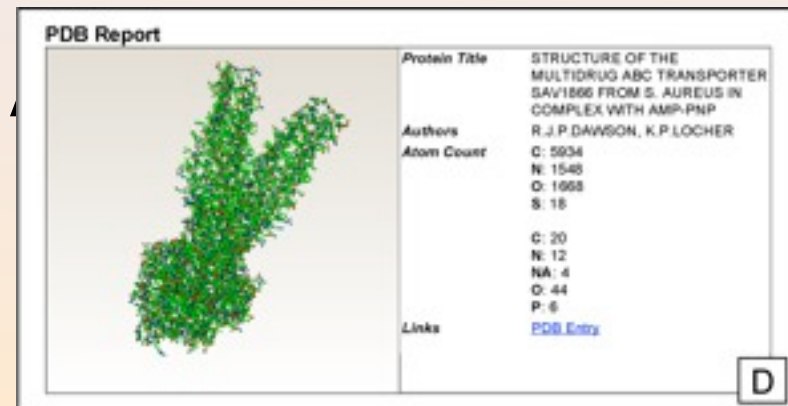   D = Δ(A,B) ∘ C may not be a valid workflow

u Find correspondences between A and C: map(A,C)
   – Diffuse similarity scores across the product graph AxC using Eigenvalue decompositions



M = MAP(A, C)

# Creating Workflows by Analogy



1. Compute difference: Δ(A,B)
   – Just like a patch!
   – But…

   D = Δ(A,B) ∘ C may not be a valid workflow

u  Find correspondences between A and C: map(A,C)
   – Diffuse similarity scores across the product graph AxC using Eigenvalue decompositions

# Creating Workflows by Analogy

1. Compute difference: Δ(A,B)
   - Just like a patch!
   - But…

   D = Δ(A,B) ∘ C may not be a valid workflow

u Find correspondences between A and C: map(A,C)
   - Diffuse similarity scores across the product graph AxC using Eigenvalue decompositions

u Compute mapped difference $\Delta_{AC}$

(A,B) =map(A,C) Δ(A,B)

A — is to — B

as — C — is to — D

1. Compute difference: Δ(A,B)
   - Just like a patch!
   - But…

   D = Δ(A,B) ∘ C may not be a valid workflow

u Find correspondences between A and C: map(A,C)
   - Diffuse similarity scores across the product graph AxC using Eigenvalue decompositions

u Compute mapped difference $\Delta_{AC}$(A,B) =map(A,C) Δ(A,B)

u D = $\Delta_{AC}$(A,B) ∘ C

# Creating Workflows by Analogy


A is to B as C is to D

1. Compute difference: Δ(A,B)
   - Just like a patch!
   - But…
   
   D = Δ(A,B) ∘ C may not be a valid workflow

u Find correspondences between A and C: map(A,C)
   - Diffuse similarity scores across the product graph AxC using Eigenvalue decompositions

u Compute mapped difference Δ(A,B) =map(A,C) Δ(A,B)

u $D = \Delta_{AC}(A,B) \circ C$

# QBE and Analogies

See paper:

- Querying and Re-Using Workflows with VisTrails
  Carlos E. Scheidegger, David Koop, Huy Vo, Juliana Freire, and Claudio T. Silva **(Best Paper Award at VIS 2007)**



Analogy Templates

Applying Analogy 1   Applying Analogy 2   Applying Analogy 3

Automatically generated workflow sequence

# VisComplete: A Workflow

- Identify graph fragments that co-occur in a collection of workflows
- Predict sets of likely workflow additions to a given partial workflow

[Koop et al., IEEE Vis 2008]



(a)

(b)

Database of Pipelines

(c)

# VisComplete: A Workflow

- Similar to a Web browser suggesting URL completions
- Idea applicable to integration queries [Sarah Cohen-Boulakia et a., JBCB 2006; Talukdar et al., VLDB 2008]

# VisComplete (video)

[Koop et al., IEEE Vis2008]

# VisComplete (video)

VisComplete:
Data-driven Suggestions for
Visualization Systems

# Acknowledgments: Funding

- This work is partially supported by the National Science Foundation, the Department of Energy, an IBM Faculty Award, and a University of Utah Seed Grant.

# More info about VisTrails

google  vistrails


Or


http://www.vistrails.org

# Emerging Work/Applications

# Visualization at large and on the go



High- resolution rendering of the Columbia River virtual estuary at a display wall

# Visualization at large and on the go



Rendering of the Columbia River on an IPOD Touch

# VisTrails: Science Dissemination

- Science mashups: simplify data exploration through visualization



[Santos et al, IEEE Vis 2009 (cond. accept)]

# Scientific Publications and Provenance



Fig. 1. Mechanical efficiency when bicycling expressed as "gross efficiency" and "delta efficiency" over the 7-yr period in this individual. WC, World Bicycle Road Racing Championships, 1st and 4th place, respectively. Tour de France 1st, Grand Champion of the Tour de France in 1999–2004.

## METHODS

*General testing sequence.* On reporting to the laboratory, training, racing, and medical histories were obtained, body weight was measured ($\pm 0.1$ kg), and the following tests were performed after informed consent was obtained, with procedures approved by the Internal Review Board of The University of Texas at Austin. Mechanical efficiency and the blood lactate threshold (LT) were determined as the subject bicycled a stationary ergometer for 25 min, with work rate increasing progressively every 5 min over a range of 50, 60, 70, 80, and 90% $Vo_{2\ max}$. After a 10- to 20-min period of active recovery, $Vo_{2\ max}$ when cycling was measured. Thereafter, body composition was determined by hydrostatic weighing and/or analysis of skin-fold thickness (34, 35).

# Scientific Publications and Provenance

"raw data from the January 1993 test that revealed several additional deviations from the *published* methodology. Coyle used a *20-min* ergometer protocol (*not 25 min*), including 2- and 3-min stages where respiratory exchange ratios (RER) exceeded 1.00. An *RER >1.00 invalidates use of the Lusk equations* (5) to estimate energy expenditure."

"…all of the published delta efficiency values are wrong. …there exists no credible evidence to support Coyle's conclusion that Armstrong's muscle efficiency improved."

*http://jap.physiology.org/cgi/content/full/105/3/1020*

# VisTrails: Science Dissemination

- Provenance-rich documents and publications

# The Provenance-Rich Paper

# Provenance and Teaching (1)

- Leverage provenance to improve the way we teach CS and Science
    - http://www.vistrails.org/index.php/SciVisFall2008
    - Lecture provenance: student can reproduce results

# Provenance and Teaching (1)

- Leverage provenance to improve the way we teach CS and Science
  - http://www.vistrails.org/index.php/SciVisFall2008
  - Lecture provenance: student can reproduce results



**Figure 5.2:** Plots of the Mauna Loa data set showing monthly measurements (left) with the yearly trend (right) using the principles for improving vision. The plot on the right is the same that was shown previously in Figure 5.1.

# Provenance and Teaching (2)

# Provenance and Teaching (2)

- Homework provenance provides insights regarding
  - Task complexity and nature: number of actions; structural vs. parameter changes; task duration
  - Student confusion: large branching factor=lots of trial and error steps
- Very detailed (and honest!) feedback: instructors can leverage this information

[Lins et al., SSDBM 2008]

# Provenance and Teaching (2)

- Homework provenance provides insights regarding
  - Task complexity and nature: number of actions; structural vs. parameter changes; task duration
  - Student confusion: large branching factor=lots of trial and error steps
- Very detailed (and honest!) feedback: instructors can leverage



Branching Structure for Task 3 of User 1

Branching Structure for Task 3 of User 2

[Lins et al., SSDBM 2008]

# Provenance and Teaching (2)

- Homework provenance provides insights regarding
  - Task complexity and nature: number of actions; structural vs. parameter changes; task duration
  - Student confusion: large branching factor=lots of trial and error steps
- Very detailed (and honest!) feedback: instructors can leverage



Branching Structure for Task 3 of User 1

[Lins et al., SSDBM 2008]

# Provenance and Teaching (3)

# Provenance and Teaching (3)

- Homework provenance helps students and instructors to *collaborate*
  - Student is stuck, sends his provenance
  - Instructor understands student's problem, provides hints---student can see what instructor did!
  - They can also collaborate in real time [Ellkvist et al., IPAW 2008]

# Using Provenance to Teach Electronic Media



[Langefeld and Kessler, Submitted 2009]

"[...] The students have gotten to the point where they demand the VisTrails files for every demonstration just after I complete [it]"

"[...] students used [a vistrail instead of a reference model] 62% of the time"

Students who used provenance produced higher-quality models

# Provenance-Based Tutorial for Maya

# Provenance Analytics: Opportunities

- ☒ Volume of collected provenance is growing
- ☒ Workflow and provenance repositories
  - – myExperiments (EU), Provenance Repository (Indiana), ManyEyes (IBM), Yahoo! Pipes
- ☒ Opportunity for knowledge discovery, sharing and re-use
  - – Discover *workflow patterns* → a recommendation system that suggests alternatives to users as they construct a workflow
  - – Discover *workflow refinement patterns* → automatically extract analogies from shared repositories
  - – *Cluster* (organize) workflow collections → simplify query and search over repositories
  - – *Infer workflow specification* from execution log [Aalst et al., TKDE 2004]

# Provenance Analytics: Challenges

- ☒ Lots of data, complex data: graphs + metadata
  - – Modules, parameters, parameter values, data products
- ☒ Do existing approaches to graph mining scale?
- ☒ Case study in clustering: [Santos, IPAW 2008]
  - – Explore different workflow representations: Graphs versus bag of words
  - – Examine trade-off between efficiency and cluster quality
  - – Bag of words surprisingly *effective*, and *much more efficient*

[NSF Medium IIS, recommended for funding, 2009]

# Conclusions and Future Work

- Advanced visualization and data analysis techniques are key to the advancement of science
- Future work into scalable algorithms, verifiable visualization, information visualization.
- Provenance management is essential for exploratory computational tasks
  - Provenance can be used to support reflective reasoning
  - Intuitive interfaces for simplifying the construction and refinement of workflows
- Science 2.0: Sharing provenance at a large scale creates new opportunities [Freire and Silva, CHI SDA, 2008]
  - Workflow/provenance repositories; provenance-enabled publications
  - Expose scientists to different techniques and tools
  - Scientists can learn by example; expedite their scientific training; and potentially reduce their time to insight
- Provenance + Workflows + Sharing have the potential to revolutionize science!

# Acknowledgments

☒ Thanks to VGC and VisTrails group

☒ This work is partially supported by the National Science Foundation, the Department of Energy, an IBM Faculty Award, and a University of Utah Seed Grant.

# Graduate Student, Postdoc, and Software Development Positions Open